

Estimation in moderately misspecified models

Nils Lid Hjort

University of Oslo and Norwegian Computing Centre

-- August 1991 --

ABSTRACT. Suppose data are fitted to some parametric model $f(y, \theta)$, but that the true model happens to be $f(y, \theta, \gamma)$, with an additional γ parameter, and where $f(y, \theta, \gamma_0) = f(y, \theta)$. When a parameter μ is to be estimated, one can use $\mu(\hat{\theta}_{\text{wide}}, \hat{\gamma}_{\text{wide}})$ based on likelihood estimation in the wider model, or $\mu(\hat{\theta}_{\text{narr}}, \gamma_0)$ based on likelihood estimation of θ alone in the narrow model. Including γ in the model means less bias but larger sampling variability. Two basic questions are addressed in this article. (i) Just how much misspecification can the narrow model tolerate? In the context of a large-sample moderate misspecification framework we find a surprisingly simple, sharp, and general answer, in the form of an explicit criterion for when narrow estimation is more precise than wide estimation. There is effectively a 'tolerance radius' around a given narrow model. The theory is illustrated by computing this tolerance radius in a selection of interesting examples, that also demonstrate the degree of robustness of important standard methods against moderate incorrectness of the model under which they are optimal. (ii) Are there other estimators that work well both under narrow and wide circumstances? We discuss several possibilities and propose some new procedures. All methods are compared in a broad performance study. This comparison can be carried out rather generally and rather simply due to a drastic reduction to a particular standard problem.

KEY WORDS: *Bayes and empirical Bayes, choice of model, compromise estimators, deliberate bias, ignorance is strength, misspecified model, parametric inference, performance study, tolerance radius*

1. Introduction and motivating examples

Our theme is moderately misspecified parametric models, and we ask two main questions. The first is: Just how much misspecification can a given parametric model tolerate in a certain direction? More specifically, when is it advantageous to stick to the narrow model, even when it is incorrect? When will 'narrow estimation' be more precise than 'wide estimation'? The second question is broader: Are there estimators that are about as good as the narrow estimator when the narrow model is correct, and about as good as the wide estimator when the narrow model is incorrect? We shall present a generous list of examples to motivate the problems and precise formulations of them.

EXAMPLE A. Suppose data Y_1, \dots, Y_n come from a life distribution on $[0, \infty)$ and that the median μ is to be estimated. If the density is the exponential $f(y) = \theta e^{-\theta y}$, then $\mu = \log 2 / \theta$, and a natural estimator is $\hat{\mu}_{\text{narr}} = \log 2 / \hat{\theta}_{\text{narr}}$, where $\hat{\theta}_{\text{narr}} = 1/\bar{Y}$ is the maximum likelihood (ML) estimator in this narrow model. If it is suspected that the model could deviate from simple exponentiality in direction of the Weibull distribution, with

$$f(y, \theta, \gamma) = \exp\{-(\theta y)^\gamma\} \gamma (\theta y)^{\gamma-1} \theta, \quad y > 0, \quad (1.1)$$

then we should conceivably use $\hat{\mu}_{\text{wide}} = (\log 2)^{1/\hat{\gamma}}/\hat{\theta}$, using ML estimators $\hat{\theta}, \hat{\gamma}$ in the wider Weibull model. But *if* the simple model is right, i.e. $\gamma = 1$, then $\hat{\mu}_{\text{narr}}$ is better, in terms (for example) of mean squared error. By sheer continuity it should be better also for γ 's close to 1. How much must γ differ from 1 in order for $\hat{\mu}_{\text{wide}}$ to become better? And what with similar questions for other typical parametric departures from exponentiality, like the gamma family?

EXAMPLE B. The most popular modelling of data Y_1, \dots, Y_n is to postulate normality, i.e. assuming $f(y) = \phi((y-\xi)/\sigma)/\sigma$ for suitable parameters ξ and σ . In many situations the normal density is too light-tailed to constitute a serious description, however. A remedy then is to use

$$f(y, \xi, \sigma, m) = g_m\left(\frac{y - \xi}{\sigma}\right) \frac{1}{\sigma},$$

where $g_m(t)$ is the t -density with m degrees of freedom. The narrower normal model corresponds to $m = \infty$, and it is naturally felt that for large m the discrepancy between normality and t -ness shouldn't matter. Suppose for example that the parameter to be estimated is sd, the standard deviation for Y_i 's. For how large m is it the case that the narrow-model estimator $\hat{\text{sd}}_{\text{narr}}$, which happens to be the ordinary empirical standard deviation, is better than the more laborious

$$\hat{\text{sd}}_{\text{wide}} = \sqrt{\frac{\hat{m}}{\hat{m} - 2}} \hat{\sigma},$$

computed from ML estimates $\hat{\xi}, \hat{\sigma}, \hat{m}$ in the three-parameter model? What with other parameters to estimate than the standard deviation?

EXAMPLE C. Consider a regression situation with n pairs (x_i, Y_i) . The classical model says $Y_i \sim N\{\alpha + \beta x_i, \sigma^2\}$ for appropriate parameters α, β, σ , and encourages for example $\hat{\mu}_{\text{narr}} = \hat{\alpha}_{\text{narr}} + \hat{\beta}_{\text{narr}} x$ as the estimator for the median (or mean value) of the distribution of Y for a given x value. Suppose however that the regression curve could be moderately quadratic, say $Y_i \sim N\{\alpha + \beta x_i + \gamma(x_i - \bar{x})^2, \sigma^2\}$ for a moderate γ . How much must γ differ from zero in order for

$$\hat{\mu}_{\text{wide}} = \hat{\alpha} + \hat{\beta}x + \hat{\gamma}(x - \bar{x})^2,$$

with regression parameters now evaluated in the wider model, to perform better? And again the same questions could be asked for other parameters, like comparing $\hat{x}_{0,\text{narr}}$ with $\hat{x}_{0,\text{wide}}$, the narrow-model based and the wide-model based estimators of the point x_0 at which the regression curve crosses a certain level.

EXAMPLE D. In some situations a more interesting departure from standard regression lies in variance heterogeneity. This could for example suggest using $Y_i \sim N\{\alpha + \beta x_i, \sigma^2(1 + \gamma x_i)\}$, where γ is zero under classical regression. For what range of γ values are standard methods, all derived under the $\gamma = 0$ hypothesis, still better than four-parameter-model analysis?

EXAMPLE E. Let us also include another type of model uncertainty, that of misspecification due to using an incorrect transformation. The transformation model invented here

has some of the intentions of the Box–Cox power transformation scheme, but avoids some of its pitfalls. It postulates that

$$h_\lambda((Y_i - \alpha - \beta x_i)/\sigma) \sim N\{0, 1\}, \quad \text{where} \quad h_\lambda(z) = \Phi^{-1}\{\Phi(z)^\lambda\}, \quad (1.2)$$

for appropriate values of $(\alpha, \beta, \sigma, \lambda)$; $\lambda = \lambda_0 = 1$ brings us back to classics. Let us briefly discuss this model and its use before we concentrate on the local misspecification part. It can be written $Y_i = \alpha + \beta x_i + \sigma Z_i$, where $h_\lambda(Z_i)$ follows the standard normal distribution for suitable transformation parameter. Varying λ gives a fair range of transformations, and in particular includes the possibility of having skewed error distributions. The four parameters can be estimated from the data. The notation is possibly deceiving, in that it invites one to think in terms of ‘ $\alpha + \beta x_i$ plus noise with level σ ’. This is not quite the case since Z_i has a skewed distribution with mean and median different from zero. It is advisable to reparameterise, after having found a suitable λ from data, to the familiar structure + noise form. One possibility is $Y_i = \{\alpha + \sigma e(\lambda)\} + \beta x_i + \sigma v(\lambda) Z_i^0$, in which $e(\lambda)$ and $v(\lambda)$ are mean value and standard deviation of Z_i , under the $h_\lambda(Z_i) \sim N\{0, 1\}$ model, and where Z_i^0 now has mean zero and standard deviation 1. Another possibility is

$$\begin{aligned} Y_i &= \{\alpha + \sigma \Phi^{-1}(0.50^{1/\lambda})\} + \beta x_i + \sigma \{\Phi^{-1}(0.75^{1/\lambda}) - \Phi^{-1}(0.25^{1/\lambda})\} Z_i' \\ &= \alpha' + \beta x_i + \sigma' Z_i', \end{aligned} \quad (1.3)$$

the point being that Z_i' has median zero and interquartile range 1. Our technical point is that (1.2) is a useful generalisation of classical regression to situations with skewed errors, and that parameter estimation is best carried out using ML machinery on (1.2); and our statistical point is that (1.3) better conveys the structure and the noise in the data, and should be used post estimation.

The present concern is how robust standard methods, which presume $\lambda = 1$, are against misspecification of that parameter. Should one use

$$\hat{\mu}_{\text{wide}}(x) = \hat{\alpha}_{\text{wide}} + \hat{\beta}_{\text{wide}}x + \hat{\sigma}_{\text{wide}}\Phi^{-1}(0.50^{1/\hat{\lambda}_{\text{wide}}}) \quad (1.4)$$

to estimate the median of Y for given x , or will the effortlessly obtainable $\hat{\mu}_{\text{narr}}(x) = \hat{\alpha}_{\text{narr}} + \hat{\beta}_{\text{narr}}x$ suffice?

EXAMPLE F. Next consider logistic regression, in which pairs (x_i, Y_i) are observed of the type $Y_i|x_i \sim \text{Bin}\{1, p(x_i)\}$, with $p(x) = \exp(\alpha + \beta x)/\{1 + \exp(\alpha + \beta x)\}$ being the standard model. Again we can ask whether standard methods based on $(\hat{\alpha}_{\text{narr}}, \hat{\beta}_{\text{narr}})$, for example for estimating the true $p(x)$ at a given x , or for estimating the cut-off point at which $p(x)$ exceeds $\frac{1}{2}$, become seriously inferior under moderate misspecifications. One natural type of departure is modelled by adding a quadratic term $\gamma(x_i - \bar{x})^2$ to the linear term; another is

$$p(x) = p(x, \alpha, \beta, \eta) = \left\{ \frac{\exp(\alpha + \beta x)}{1 + \exp(\alpha + \beta x)} \right\}^\eta, \quad (1.5)$$

where it is of interest to vary η around $\eta_0 = 1$.

EXAMPLE G. Our final example is the two-sample model with variances that may or may not be equal. So X_1, \dots, X_m are $N\{\xi_1, \sigma_1^2\}$ and Y_1, \dots, Y_n are $N\{\xi_2, \sigma_2^2\}$, all of them

are independent, and the narrow model specifies that $\sigma_1 = \sigma_2$. Under this assumption it is easy to put up estimators, confidence intervals etc. for parameters related to the difference between the X -distribution and the Y -distribution, like the Mahalanobis distance $\Delta = |\xi_2 - \xi_1|/\sigma$. More awkward methods are needed when $\sigma_2 \neq \sigma_1$, cf. the Behrens-Fisher problem. The in some sense natural generalisation of the Mahalanobis distance is for example

$$\Delta = (\nu^2 + \omega^2)^{1/2}, \quad \text{where } \nu^2 = (\xi_2 - \xi_1)^2/\sigma^2, \quad \omega^2 = 4 \log \frac{\sigma^2}{\sigma_1 \sigma_2}, \quad \sigma^2 = (\sigma_1^2 + \sigma_2^2)/2,$$

see Hjort (1986a, Ch. 10). How resistant is the simple $\hat{\Delta}_{\text{narr}} = |\bar{Y} - \bar{X}|/\hat{\sigma}_{\text{narr}}$ to differences in σ_1, σ_2 ? When is it necessary to use the much more complicated $\hat{\Delta}_{\text{wide}}$?

Let us summarise the common characteristics of these situations. There is a narrow and usually simple parametric model which can be fitted to the data, but there is a potential misspecification, which can be ameliorated by its encapsulation in a wider model with one additional parameter. Estimating a parameter assuming correctness of the narrow model involves modelling bias, but doing it in the wider model could mean larger sampling variability. Thus the choice of method becomes a statistical balancing act with perhaps deliberate bias against variance.

The examples above span a reasonable range of heavily used 'narrow' models along with indications of rather typical kinds of deviances from them. Many standard textbook methods for parametric inference are derived under the conditions of such narrow models. Our main result, derived in Section 3, is a surprisingly sharp and general large-sample criterion for how much misspecification a given narrow model can tolerate. This criterion is applied to Examples A-G in Section 7. It is relatively easy to compute, in that it only involves the familiar Fisher information matrix, for the wide model, but evaluated under narrow model conditions. A particularly pleasing facet of our tolerance criterion is that it does not depend upon the particular parameter estimand at all!

In addition to quantifying the degree of robustness of standard methods there are also pragmatic reasons for the present investigation. Statistical analysis will in practice still be carried out using narrow model based methods in the majority of cases, for reasons of ignorance, simplicity, naïveté and boldness; using wide model methods will very often be much more laborious, and only experts will use them anyhow. Thus it is of interest to quantify the consequences of ignorance, and it would be nice to obtain permission to go on doing analysis as if the simple model were true. Such a partial permission is in fact given here. The results of this paper can be interpreted as saying that mild departures from the narrow model do not really matter, and that in fact more ambitious methods could perform worse. In the examples of Section 7 quite explicit limits are given for the degree of misspecification that is tolerable. But this upper limit is in most cases dependent upon parameters of the model, and should be estimated by the conscientious statistician in situations where departures of the type described are suspected. So one should by all means carry out estimation of the additional parameter, even if it turns out that it won't be needed in the final analysis.

Several tangential topics are taken up in Section 4. These include measures of distance from null model to the least tolerable misspecification; comparison with the model selection

criteria of Akaike and Schwarz; simulation based evaluation of our criterion; discussion of the concept of a robust model; dangerous versus noncritical departures from a model; interpretation of confidence intervals under misspecification; and deviances from a model in more than one direction.

There is also room for improvement over the narrow and wide methods. In Section 5 some new estimators are proposed that are designed to work well both under narrow and wide circumstances. A broad comparison of the various compromise estimators is made, in a large-sample framework of moderately misspecified parametric models. A connection to Bayesian robustness is also made. We are able to make a quite general and drastic reduction: The performance of a large class of competing estimators can be studied in a much simpler and very classical context, that of estimating a in a $N\{a, 1\}$ situation with one observation! Here the narrow model corresponds to $a = 0$. This provides fresh motivation for studying a -estimators that in various ways take into account that values of a in the vicinity of zero are perhaps more likely or perhaps more important. Such a study is reported on in Sections 5 and 6.

The traditional robustness literature is mostly concerned with construction of methods that perform well over a ‘nonparametric neighbourhood’ around some basic model. The present work is different in that it envisages specific, parametric alternatives to the basic model. There is a literature on parametric robustness, perhaps chiefly concerned with studying behaviour of standard methods and modified standard methods under natural violations of the basic model. Only rarely have comparisons been made between ‘narrow’ and ‘wide’ methods, however. Some papers have calculated and commented on the increased estimation noise for a narrow model parameter when passing to a wider model, like comparing the variances of $\hat{\theta}_{\text{narr}}$ and $\hat{\theta}_{\text{wide}}$ in Example A. This is beside the point, partly confusing, and not very interesting, since what matters is studying ‘real’ parameters which are meaningful functions of the full model, as the median $\mu = \mu(f) = \mu(\theta, \gamma) = (\log 2)^{1/\gamma}/\theta$ in Example A, the standard deviation $\text{sd} = \text{sd}(f) = \text{st}(\xi, \sigma, m) = \{m/(m-2)\}^{1/2}\sigma$ in Example B, etcetera. Bickel (1984) is on the other hand clear about this issue, and is concerned with several problems that resemble those considered here. He does not compare narrow and wide methods, and does not study tolerance distances, but works directly with certain minimax strategies, in a framework of nested linear normal models; see also 5G below. The paper by Berger (1982) on Bayesian robustness also turns out to be related to some of these questions. See Bickel’s comments on Berger and 5E, 5F, 5G below.

2. Large-sample framework for the problem

We shall start our investigation in the i.i.d. framework. Suppose Y_1, \dots, Y_n come from some common density f , and represent the wide model as $f(y) = f(y, \theta, \gamma)$, where $\gamma = \gamma_0$ corresponds to the narrow model, say $f(y, \theta) = f(y, \theta, \gamma_0)$. We assume that $\theta = (\theta_1, \dots, \theta_p)'$ lies in some open region in Euclidian p -space, that γ lies in some interval containing γ_0 , and that the wide model is ‘smooth’; for definiteness we postulate that the regularity conditions put forward in Lehmann’s (1983) Chapter 6.4 are in force. We are to study behaviour of estimators when γ deviates from γ_0 . The parameter to be estimated is some $\mu = \mu(f)$, which we write as $\mu(\theta, \gamma)$ since the wider model is assumed to be an adequate description of reality. We concentrate on ML procedures, and write $\hat{\theta}_{\text{narr}}$ for the

estimator of θ in the narrow model and $(\hat{\theta}, \hat{\gamma})$ for the estimators in the wide model. The two major entries in the competition are

$$\hat{\mu}_{\text{narr}} = \mu(\hat{\theta}_{\text{narr}}, \gamma_0) \quad \text{and} \quad \hat{\mu}_{\text{wide}} = \mu(\hat{\theta}, \hat{\gamma}) \quad (2.1)$$

(but see Section 5 for other estimators).

These could be compared in an asymptotic framework in which Y_i 's come from some fixed $f(y, \theta_0, \gamma)$, and $\gamma \neq \gamma_0$. In this case $\sqrt{n}(\hat{\mu}_{\text{wide}} - \mu)$ has a limit distribution, which can be derived from the proposition below. The situation is different for the narrow model procedure. Here $\sqrt{n}(\hat{\mu}_{\text{narr}} - \mu)$ can be represented as a sum of two terms. The first is $\sqrt{n}\{\mu(\hat{\theta}_{\text{narr}}, \gamma_0) - \mu(\theta_0, \gamma_0)\}$, which has a limit distribution, with generally smaller variability than that of the wide model procedure, and the second is $-\sqrt{n}\{\mu(\theta_0, \gamma) - \mu(\theta_0, \gamma_0)\}$, which tends to plus or minus infinity, reflecting a bias that for very large n will dominate completely. This merely goes to show that with very large sample sizes one is penalised for any bias and one should use the wide model. This result is somewhat irrelevant, however, and suggests that a large-sample framework which uses a local neighbourhood of γ_0 that shrinks when the sample size grows is much more adequate. Study therefore model P_n , the n 'th model, under which

$$Y_1, \dots, Y_n \text{ are i.i.d. from } f_n(y) = f(y, \theta_0, \gamma_0 + \delta/\sqrt{n}), \quad (2.2)$$

and where θ_0 is fixed but arbitrary. In this framework we need limit distributions for the wide model estimators $(\hat{\theta}, \hat{\gamma})$ and for the narrow model estimator $\hat{\theta}_{\text{narr}}$.

Consider

$$\begin{pmatrix} U(y) \\ V(y) \end{pmatrix} = \begin{pmatrix} \partial \log f(y, \theta_0, \gamma_0) / \partial \theta \\ \partial \log f(y, \theta_0, \gamma_0) / \partial \gamma \end{pmatrix}, \quad (2.3)$$

the score function for the wide model, but evaluated at the null point (θ_0, γ_0) ; and the accompanying familiar $(p+1) \times (p+1)$ size information matrix

$$J_{\text{wide}} = \text{VAR}_0 \begin{pmatrix} \partial \log f(Y, \theta_0, \gamma_0) / \partial \theta \\ \partial \log f(Y, \theta_0, \gamma_0) / \partial \gamma \end{pmatrix} = \begin{pmatrix} J_{11} & J_{12} \\ J_{21} & J_{22} \end{pmatrix}.$$

Note that the $p \times p$ size J_{11} is simply the information matrix of the narrow model, evaluated at θ_0 , and that the scalar J_{22} is the variance of $V(Y_i)$, also computed under the narrow model.

PROPOSITION. *Under the sequence of models P_n of (2.2), as n tends to infinity, we have*

$$\begin{pmatrix} \sqrt{n}(\hat{\theta} - \theta_0) \\ \sqrt{n}(\hat{\gamma} - \gamma_0 - \delta/\sqrt{n}) \end{pmatrix} \rightarrow_d N_{p+1}\{0, J_{\text{wide}}^{-1}\}, \text{ or } \begin{pmatrix} \sqrt{n}(\hat{\theta} - \theta_0) \\ \sqrt{n}(\hat{\gamma} - \gamma_0) \end{pmatrix} \rightarrow_d N_{p+1}\left\{\begin{pmatrix} 0 \\ \delta \end{pmatrix}, J_{\text{wide}}^{-1}\right\};$$

$$\sqrt{n}\{\hat{\theta}_{\text{narr}} - (\theta_0 + J_{11}^{-1} J_{12} \delta / \sqrt{n})\} \rightarrow_d N_p\{0, J_{11}^{-1}\}, \text{ or } \sqrt{n}(\hat{\theta}_{\text{narr}} - \theta_0) \rightarrow_d N_p\{J_{11}^{-1} J_{12} \delta, J_{11}^{-1}\}.$$

PROOF: Consider $\hat{\theta}_{\text{narr}}$ first. The familiar Taylor expansion arguments that lead to the classical $\sqrt{n}(\hat{\theta}_{\text{narr}} - \theta_0) \rightarrow_d N\{0, J_{11}^{-1}\}$ under the null model $f(x, \theta_0, \gamma_0)$ can be used in the present $\gamma_0 + \delta/\sqrt{n}$ case as well. For

$$\sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(Y_i, \hat{\theta}_{\text{narr}}, \gamma_0) = \sum_{i=1}^n U(Y_i) + I_n(\tilde{\theta}_n)(\hat{\theta}_{\text{narr}} - \theta_0) = 0,$$

in which $I_n(\theta) = \sum_{i=1}^n \frac{\partial^2}{\partial \theta \partial \theta} \log f(Y_i, \theta, \gamma_0)$ and $\tilde{\theta}_n$ lies somewhere between θ_0 and $\hat{\theta}_{\text{narr}}$. Under the conditions stated $\hat{\theta}_{\text{narr}} \rightarrow_p \theta_0$, under P_n , using necessary but moderate variations of the arguments used in Lehmann's (1983) Chapter 6.4 and 6.8, and $-I_n(\theta_0)/n$ as well as $-I_n(\tilde{\theta}_n)/n$ tend in probability, under P_n , to J_{11} . All this leads to

$$\sqrt{n}(\hat{\theta}_{\text{narr}} - \theta_0) \doteq_d \{-\frac{1}{n}I_n(\theta_0)\}^{-1} \sqrt{n}\bar{U}_n \doteq_d J_{11}^{-1} \sqrt{n}\bar{U}_n, \quad (2.4)$$

where $A_n \doteq_d B_n$ means that $A_n - B_n$ tends to zero in probability, and \bar{U}_n is the average of the n first $U(Y_i)$'s. The triangular version of the Lindeberg theorem shows that $\sqrt{n}\bar{U}_n$ tends in distribution, under P_n , to $N_p\{J_{12}\delta, J_{11}\}$. This is because

$$\begin{aligned} E_{P_n} U(Y_i) &= \int f(y, \theta_0, \gamma_0 + \delta/\sqrt{n}) U(y) dy \\ &\doteq \int f(y, \theta_0, \gamma_0) \{1 + V(y)\delta/\sqrt{n}\} U(y) dy = J_{12}\delta/\sqrt{n}, \end{aligned}$$

and similarly $U(Y_i)U(Y_i)'$ can be shown to have expected value $J_{11} + O(\delta/\sqrt{n})$, under P_n . This proves the 'narrow' part of the proposition.

Similar reasoning takes care of the 'wide' part too. One finds

$$\begin{pmatrix} \sqrt{n}(\hat{\theta} - \theta_0) \\ \sqrt{n}(\hat{\gamma} - \gamma_0) \end{pmatrix} \doteq_d J_{\text{wide}}^{-1} \begin{pmatrix} \sqrt{n}\bar{U}_n \\ \sqrt{n}\bar{V}_n \end{pmatrix} \rightarrow_d J_{\text{wide}}^{-1} N_{p+1}\left\{ \begin{pmatrix} J_{12}\delta \\ J_{22}\delta \end{pmatrix}, J_{\text{wide}} \right\}, \quad (2.5)$$

which is equivalent to the wide part statement. \square

REMARK. Let us for a moment consider more general departures from the $f(y, \theta)$ model. Assume only that data Y_i come from a fixed f . Then one can show that $\hat{\theta}_{\text{narr}}$ is consistent for the particular 'least false' or 'most appropriate' parameter value $\theta_{1.f.} = \theta(f)$ that minimises the Kullback-Leibler distance $d[f: f(\cdot, \theta)] = \int f(y) \log\{f(y)/f(y, \theta)\} dy$. One can also show that $\sqrt{n}(\hat{\theta}_{\text{narr}} - \theta_{1.f.})$ tends in distribution to $N_p\{0, J(f)^{-1}K(f)J(f)^{-1}\}$, in which

$$J(f) = -E_f \frac{\partial^2 \log f(Y, \theta_{1.f.})}{\partial \theta \partial \theta}, \quad \text{and} \quad K(f) = \text{VAR}_f \frac{\partial \log f(Y, \theta_{1.f.})}{\partial \theta}. \quad (2.6)$$

This is for example made clear in Hjort (1986a, Ch. 3). — Let us apply this to the local misspecification situation, that is, insert $f(y) = f(y, \theta_0, \gamma)$, where γ is close to γ_0 . Then, by judicious Taylor expansion arguments, one can show that

$$\theta_{1.f.} = \theta_0 + J_{11}^{-1} J_{12}(\gamma - \gamma_0) + O((\gamma - \gamma_0)^2), \quad J(f)^{-1} K(f) J(f)^{-1} = J_{11}^{-1} + O(\gamma - \gamma_0).$$

Using this, for the local $\gamma = \gamma_0 + \delta/\sqrt{n}$, can be used to prove the ‘narrow part’ of the proposition again. Note that the notion and interpretation of a best fitting parameter changes when the model changes, and that the results about $\theta_{\text{l.f.}}$ quantify this in a precise way.

3. Solution

In the large-sample framework of the previous section we are to compare two estimators: The ‘safe’ $\hat{\mu}_{\text{wide}} = \mu(\hat{\theta}, \hat{\gamma})$ based on ML estimation in the wide model, and the ‘risky’ $\hat{\mu}_{\text{narr}} = \mu(\hat{\theta}_{\text{narr}}, \gamma_0)$. The true parameter is $\mu_{\text{true}} = \mu(\theta_0, \gamma_0 + \delta/\sqrt{n})$ under P_n , the n ’th model.

First consider the safe estimator. By the delta method of linearisation we find

$$\begin{aligned} & \sqrt{n}\{\mu(\hat{\theta}, \hat{\gamma}) - \mu(\theta_0, \gamma_0 + \delta/\sqrt{n})\} \\ & \doteq_d \left(\frac{\partial \mu}{\partial \theta}\right)' \sqrt{n}(\hat{\theta} - \theta_0) + \left\{\left(\frac{\partial \mu}{\partial \gamma}\right) + O(1/\sqrt{n})\right\} \sqrt{n}(\hat{\gamma} - (\gamma_0 + \delta/\sqrt{n})) \rightarrow_d N\{0, \tau^2\}, \end{aligned}$$

where

$$\tau^2 = \begin{pmatrix} \frac{\partial \mu}{\partial \theta} \\ \frac{\partial \mu}{\partial \gamma} \end{pmatrix}' J_{\text{wide}}^{-1} \begin{pmatrix} \frac{\partial \mu}{\partial \theta} \\ \frac{\partial \mu}{\partial \gamma} \end{pmatrix}. \quad (3.1)$$

The partial derivatives are computed at the null point (θ_0, γ_0) . Similarly, for the risky estimator,

$$\begin{aligned} & \sqrt{n}\{\mu(\hat{\theta}_{\text{narr}}, \gamma_0) - \mu(\theta_0, \gamma_0 + \delta/\sqrt{n})\} \\ & = \sqrt{n}\{\mu(\hat{\theta}_{\text{narr}}, \gamma_0) - \mu(\theta_0, \gamma_0)\} - \sqrt{n}\{\mu(\theta_0, \gamma_0 + \delta/\sqrt{n}) - \mu(\theta_0, \gamma_0)\} \\ & \doteq_d \left(\frac{\partial \mu}{\partial \theta}\right)' \sqrt{n}(\hat{\theta}_{\text{narr}} - \theta_0) - \sqrt{n} \frac{\partial \mu}{\partial \gamma} \delta / \sqrt{n} \rightarrow_d N\{b\delta, \tau_0^2\}, \end{aligned}$$

in which

$$b = J_{21} J_{11}^{-1} \left(\frac{\partial \mu}{\partial \theta}\right) - \frac{\partial \mu}{\partial \gamma} \quad \text{and} \quad \tau_0^2 = \left(\frac{\partial \mu}{\partial \theta}\right)' J_{11}^{-1} \left(\frac{\partial \mu}{\partial \theta}\right). \quad (3.2)$$

By evaluating the mean value of the square of the limit distributions we have that n times the asymptotic mean squared error of $\hat{\mu}_{\text{wide}}$ becomes τ^2 , while the corresponding quantity for $\hat{\mu}_{\text{narr}}$ becomes $b^2 \delta^2 + \tau_0^2$.

We are now in a position to find out when the risky estimator is better than the safe one, simply by algebraically solving the inequality $b^2 \delta^2 + \tau_0^2 \leq \tau^2$ w.r.t. δ . Start out writing

$$J_{\text{wide}}^{-1} = \begin{pmatrix} J^{11} & J^{12} \\ J^{21} & J^{22} \end{pmatrix},$$

where a prominent rôle is designated for

$$J^{22} = \kappa^2 = (J_{22} - J_{21} J_{11}^{-1} J_{12})^{-1} \quad (3.3)$$

in what follows, and $J^{12} = -J_{11}^{-1} J_{12} \kappa^2$, $J^{11} = J_{11}^{-1} + J_{11}^{-1} J_{12} J_{21} J_{11}^{-1} \kappa^2$. This leads to the simplification

$$\begin{aligned} \tau^2 & = \left(\frac{\partial \mu}{\partial \theta}\right)' J_{11}^{-1} \left(\frac{\partial \mu}{\partial \theta}\right) + \left(\frac{\partial \mu}{\partial \theta}\right)' J_{11}^{-1} J_{12} J_{12}' J_{11}^{-1} \left(\frac{\partial \mu}{\partial \theta}\right) \kappa^2 - 2 \left(\frac{\partial \mu}{\partial \theta}\right)' J_{11}^{-1} J_{12} \left(\frac{\partial \mu}{\partial \gamma}\right) \kappa^2 + \left(\frac{\partial \mu}{\partial \gamma}\right)^2 \kappa^2 \\ & = \tau_0^2 + b^2 \kappa^2. \end{aligned}$$

We have reached

RESULT. (i) The case where $b = 0$ is rather trivial; this typically corresponds to asymptotic independence between $\hat{\theta}$ and $\hat{\gamma}$ under the null model, and a parameter μ functionally independent of γ . In this case $\hat{\mu}_{\text{wide}}$ and $\hat{\mu}_{\text{narr}}$ are asymptotically equivalent, regardless of δ . (ii) In the more interesting case $b \neq 0$, the narrow model based estimator is better than or as good as the wider model based estimator if and only if $\delta^2 \leq \kappa^2$, or $|\delta| \leq \kappa$, or $|\gamma - \gamma_0| \leq \kappa/\sqrt{n}$.

Extension to regression models. To solve the problems raised in the regression type examples of the introduction we also need the similar result in the more general situation of independent observations with covariates. This can be done in a fairly straightforward fashion. Examples C–F of Sections 1 and 7 lead us naturally to the following general framework. Suppose (x_i, Y_i) are independent pairs, where Y_i has density $f(y_i, \sigma, \beta, \gamma | x_i)$ for given x_i -value, carrying some scale parameter σ (but not necessarily), a vector $\beta = (\beta_1, \dots, \beta_p)'$ of ordinary regression parameters, plus some interesting one-dimensional extra parameter γ that signals departure from the underlying classical model, which corresponds to some appropriate $\gamma = \gamma_0$.

Under mild regularity conditions the main result above continues to be true for regression models, with κ^2 defined as in (3.3), but with a somewhat more cumbersome J_{wide} matrix than before. The correct definition is now

$$J_{\text{wide}} = \lim_{n \rightarrow \infty} J_{n, \text{wide}} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \text{VAR}_0 \begin{pmatrix} \partial \log f(Y_i, \sigma_0, \beta_0, \gamma_0 | x_i) / \partial \sigma \\ \partial \log f(Y_i, \sigma_0, \beta_0, \gamma_0 | x_i) / \partial \beta \\ \partial \log f(Y_i, \sigma_0, \beta_0, \gamma_0 | x_i) / \partial \gamma \end{pmatrix}, \quad (3.4)$$

where the variance matrices are computed at the null model, under $(\sigma_0, \beta_0, \gamma_0)$. The necessary regularity conditions can be put up in various forms. These would be Lindebergian to secure normal limits and must in particular imply convergence of $J_{n, \text{wide}}$; this usually follows if it is assumed that the collection of x_i 's come from some distribution in the design space. In practice one would typically use $J_{n, \text{wide}}$ to compute κ^2 . Examples are given in Section 7.

4. Discussion

4A. Simplicity. It is remarkable that the criterion reached does not depend on the particularities of the specific parameter $\mu(\theta, \gamma)$ at all. Thus, in the situation of Example A in the introduction, calculations in Section 7 show that $|\gamma - 1| \leq 1.245/\sqrt{n}$ guarantees that being simple-minded, assuming exponentiality, works better than being ambitious, using a gamma-family, for *every* smooth parameter $\mu(\theta, \gamma)$.

Our criterion $\delta^2 \leq \kappa^2$ can be evaluated and assessed just from knowledge of J_{wide} , the information matrix of the full model, but computed at the narrow model only. This is fortunate, as the general $p + 1$ parameter matrix will be very hard to compute in many applications, but will be simpler and manageable at the null model. This is demonstrated in Section 7. Observe that the $|\delta| \leq \kappa$ criterion can be thought of in terms of the limiting variance for $\hat{\gamma}$, at the null model, since $\sqrt{n}(\hat{\gamma} - \gamma_0)$ tends to $N\{0, \kappa^2\}$ then.

4B. *How far away is the border line?* We have shown that the simple θ parameter model can tolerate up to $\gamma_0 + \kappa/\sqrt{n}$ deviation from γ_0 in the encapsulating (θ, γ) model. How far is the border line $\delta = \kappa$ from the narrow model? One way of answering this is in terms of the probability of actually detecting that the narrow model is wrong. The natural 5% level test for the correctness of the narrow model, against the alternative hypothesis that the additional γ parameter must be included, is to reject when $Z_n^2 = n(\hat{\gamma} - \gamma_0)^2/\hat{\kappa}^2$ exceeds 1.96^2 , since Z_n^2 has a limiting χ_1^2 distribution under the narrow model. Here $\hat{\kappa}$ is any consistent estimator of κ , or simply equal to the known value in such cases. The probability that this test detects that γ is not equal to γ_0 , when it in fact is equal to $\gamma_0 + \delta/\sqrt{n}$, converges to

$$\text{power}(\delta) = \Pr\{\chi_1^2(\delta^2/\kappa^2) > 1.96^2\}, \quad (4.1)$$

featuring the non-central chi squared with 1 degree of freedom and eccentricity parameter δ^2/κ^2 . This is a consequence of the proposition proved in Section 2. In particular the approximate power at the border case is equal to 17.0%. We can therefore restate the basic result as follows: Provided the true model deviates so modestly from the narrow model, that the probability of detecting it is 17.0% or less with the natural 5% level test, then the risky estimator is better than the safe estimator. Corresponding other figures for (level, power) are, for illustration, (0.01, 0.057), (0.10, 0.264), (0.20, 0.400), (0.29, 0.500).

4C. *Other distance measures.* Let us present a couple of further measures of the distance from null model to border line misspecification. (i) The Kullback–Leibler distance $d[f(\cdot, \theta_0, \gamma_0): f(\cdot, \theta_0, \gamma_0 + \delta/\sqrt{n})]$ can by clever Taylor expansion arguments be shown to be equal to $\frac{1}{2}\delta^2 J_{22}/n$ plus smaller terms, and in the border case the distance becomes $\kappa^2 J_{22}/2n$. (ii) Next consider the so-called statistical distance or L_1 -distance between the two neighbouring distributions. It is

$$\int |f(y, \theta_0, \gamma_0 + \delta/\sqrt{n}) - f(y, \theta_0, \gamma_0)| dy \doteq \frac{\delta}{\sqrt{n}} \int |V(y)| f(y, \theta_0, \gamma_0) dy.$$

This distance has a direct probabilistical interpretation. In Example A, for example, the L_1 -distance from exponentiality to the least tolerable Weibull, becomes about $0.923/\sqrt{n}$. (iii) Finally consider weighted L_2 -distance $\int (f - f_0)^2/f_0 dy$. An approximation is seen to be $\delta^2 J_{22}/n$, and the least tolerable distance is $\kappa^2 J_{22}/n$. — Note that these three distance measures are transformation invariant. See also 4F.

4D. *Comparison with Akaike's Information Criterion.* The misspecification problem is related to that of choosing a model. One general method for doing this is to use the information criterion of Akaike. In the present setting one is to compare

$$\text{AIC}_{\text{narr}} = \log L_{\text{max, narr}} - p \quad \text{with} \quad \text{AIC}_{\text{wide}} = \log L_{\text{max, wide}} - (p + 1),$$

featuring maximised log likelihoods under respectively the narrow model with p parameters and the wide model with $p + 1$ parameters. The method consists of choosing the model with largest observed AIC.

It is instructive to study AIC's behaviour in the framework of this article. Using Taylor expansion, along with techniques and notation as in the proof of the proposition of Section 2, one can show that

$$\begin{aligned} \text{AIC}_{\text{narr}} &\doteq_d \sum_{i=1}^n \log f(Y_i, \theta_0, \gamma_0) + n \bar{U}_n' J_{11}^{-1} \bar{U}_n - p, \\ \text{AIC}_{\text{wide}} &\doteq_d \sum_{i=1}^n \log f(Y_i, \theta_0, \gamma_0) + n \begin{pmatrix} \bar{U}_n \\ \bar{V}_n \end{pmatrix}' J_{\text{wide}}^{-1} \begin{pmatrix} \bar{U}_n \\ \bar{V}_n \end{pmatrix} - (p+1). \end{aligned}$$

Consequently

$$\begin{aligned} \text{AIC}_{\text{wide}} - \text{AIC}_{\text{narr}} &\doteq_d n \begin{pmatrix} \bar{U}_n \\ \bar{V}_n \end{pmatrix}' J_{\text{wide}}^{-1} \begin{pmatrix} \bar{U}_n \\ \bar{V}_n \end{pmatrix} - n \bar{U}_n' J_{11}^{-1} \bar{U}_n - 1 \\ &= n [\bar{U}_n' (J^{11} - J_{11}^{-1}) \bar{U}_n + 2 \bar{U}_n' J^{12} \bar{V}_n + \bar{V}_n' J^{22} \bar{V}_n] - 1 \\ &= n (\bar{V}_n - \bar{U}_n' J_{11}^{-1} J_{12})^2 \kappa^2 - 1 \\ &\rightarrow_d N\{-\delta/\kappa^2, 1/\kappa^2\}^2 \kappa^2 - 1 = \chi_1^2(\delta^2/\kappa^2) - 1. \end{aligned} \quad (4.2)$$

The probability that AIC prefers the wide model over the narrow model is therefore approximately $\Pr\{\chi_1^2(\delta^2/\kappa^2) > 1\}$. In particular, if the narrow model is perfect, the probability is 0.317, and in the border-line case suggested by this article, i.e. $\delta = \kappa$, the probability is 0.523. All in all the Akaike method agrees very well with the implicit advice of Section 3.

It is also instructive to see that $\text{AIC}_{\text{wide}} - \text{AIC}_{\text{narr}}$ above is asymptotically equivalent to $Z_n^2 - 1$, where $Z_n = \sqrt{n}(\hat{\gamma} - \gamma_0)/\hat{\kappa} \rightarrow_d N\{\delta/\kappa, 1\}$. This is the test statistic also discussed in 4B. Akaike's method is also related to the pre-test type strategy discussed in Section 5.

Akaike's criterion has a reputation for overfitting too often, and researchers often use a more stingy criterion due to Schwarz. It penalises with the factor $\frac{1}{2} \log n$ times the number of parameters in the model, i.e. subtracts $(\frac{1}{2} \log n)p$ and $(\frac{1}{2} \log n)(p+1)$ instead. The reasoning above, applied to this alternative criterion, shows that the Schwarz method chooses the narrow model, with probability tending (but slowly) to 1. The alternative model must be at least $\delta(\log n)^{1/2}/\sqrt{n}$ away to interest Schwarz [SIC].

4E. Evaluation of κ through stochastic simulation. The examples of Section 7 show that it is possible to compute J_{wide} and κ^2 explicitly even for somewhat complicated departure models, in effect because the computations only need to be carried out at the null model. In some situations it might be too difficult, however. One way out is then to write down the difficult elements of the J_{wide} matrix in terms of integrals, involving the null density $f(y, \theta_0, \gamma_0)$ as well as $U(y)$ and $V(y)$, and then carry out numerical integration. This is feasible since only one-dimensional integrals are involved. This method gives a numerical value of κ for specified basis point θ_0 .

Another way is through stochastic simulation. Several options can be considered. (i) Simulate a large number of Y_i 's from the null distribution, and compute score functions $U(Y_i)$ and $V(Y_i)$ along the way (see (2.3)). Then compute empirical covariances and variances to get J_{wide} . (ii) Keep n fixed, simulate Y_1^*, \dots, Y_n^* from the null density, at some desired θ_0 , and compute the estimates $\hat{\theta}^*$ and $\hat{\gamma}^*$ based on this pseudo-sample. Do

this a large number of times, and the empirical covariance matrix for $(\hat{\theta}^*, \hat{\gamma}^*)$ is J_{wide}^{-1}/n .
 (iii) Or drop $\hat{\theta}^*$ and just evaluate the empirical standard deviation of $\sqrt{n}(\hat{\gamma}^* - \gamma_0)$, which is κ . This is a feasible approach in complex regression models, or in parametric and semiparametric survival data models with censoring, where analytical expressions for κ^2 cannot be found.

4F. Good models and dangerous departures. Which departures from a given narrow model are dangerous, and which are insignificant? And what qualities should a ‘good and robust’ model have?

We have demonstrated that the narrow model can tolerate $\delta = \sqrt{n}(\gamma - \gamma_0)$ up to the limit κ in absolute value. The numerical value of κ depends on the scale used, however. The appropriate scale invariant tolerance measure is $d = \kappa^2 J_{22} = J^{22} J_{22}$, as is also suggested by the distances considered in 4C. Two numbers of this kind can be directly compared for two specifically envisaged model departures. A model departure with large d is less dangerous than one with small d .

A model deviance can be studied in terms of $V(y) = \partial \log f(y, \theta_0, \gamma_0) / \partial \gamma$, see (2.3). How well is $V(y)$ explained by the existing model, represented by $U(y)$? A natural measure is the so-called maximal correlation, $\rho^2\{U, V\}$, the maximal value of $\text{corr}\{a_1 U_1 + \dots + a_p U_p, V\}^2$ as $a = (a_1, \dots, a_p)'$ varies. It is well known and just a piece of linear algebra to prove that $a_0 = J_{11}^{-1} J_{12}$ maximises, with resulting

$$\rho^2\{U, V\} = J'_{12} J_{11}^{-1} J_{12} = 1 - 1/(\kappa^2 J_{22}) = 1 - 1/d. \quad (4.3)$$

This invites a geometrical interpretation for the tolerance limit d . The smallest possible value for d is 1, which happens when the model departure is ‘completely new’ and orthogonal to the existing model, with $J_{12} = 0$. Only a mild departure in this direction can be tolerated. So a dangerous departure is one that can be realistically suspected, in the first place, and which has a small d , or a small correlation. A non-critical departure is one that has a large tolerance d , or a large correlation, or one that perhaps is unrealistic a priori. — A good and robust model, therefore, is one where realistically suspected deviances have large tolerances d . See the examples of Section 7.

4G. Can we de-bias? We have demonstrated that narrow estimation, which means introducing a deliberate bias to reduce variability, leads to better estimator precision in a certain radius around the narrow model. The precise quantitative result is that $\sqrt{n}(\hat{\mu}_{\text{narr}} - \mu_{\text{true}})$ tends to $N\{b\delta, \tau_0^2\}$, see Section 3. Can we remove the bias and do even better?

About the best we can do in this direction is to use $\hat{\mu}_{\text{db}} = \hat{\mu}_{\text{narr}} - b(\hat{\gamma} - \gamma_0)$. Analysis reveals, working from the basis result (5.2) of the next section, that $\sqrt{n}(\hat{\mu}_{\text{db}} - \mu_{\text{true}})$ tends to $N\{0, \tau_0^2 + b^2 \kappa^2\}$. So the bias can be removed, but the price one pays amounts exactly to what was won by deliberate biasing in the first place, and the de-biased estimator is equivalent to $\hat{\mu}_{\text{wide}}$. The reason for the extra variability is that no consistent estimator exists for δ .

4H. Dwindling confidence. We have established that $\hat{\mu}_{\text{narr}}$ has higher precision than $\hat{\mu}_{\text{wide}}$ for moderate misspecifications of the narrow model. But what with further inference?

Consider confidence intervals. The usual approximate 90% interval for μ based on narrow model assumptions is $\text{CI}_{\text{narr}} = \hat{\mu}_{\text{narr}} \pm 1.645 \hat{\tau}_0 / \sqrt{n}$, where $\hat{\tau}_0$ is consistent for τ_0

of (3.2). But in the present local misspecification framework $\sqrt{n}(\hat{\mu}_{\text{narr}} - \mu_{\text{true}})$ tends to $N\{b\delta, \tau_0^2\}$, and the bias destroys the 90% property. The probability that μ_{true} is covered by CI_{narr} converges to $\Pr[-1.645 \leq N\{b\delta/\tau_0, 1\} \leq 1.645]$. This is always *strictly less than 90%*, unless the narrow model is exactly true or b of (3.2) is zero. Yes, I am shocked. The difference is not necessarily dramatic, in that the coverage probability is above 85% when $|b\delta|/\tau_0$ is smaller than 0.54 and above 80% when the ratio is smaller than 0.77. What is important is that the narrow model based interval always underestimates the confidence, under any model departure from any given parametric model, and that we have an illuminating explicit formula for the true (asymptotic) coverage probability.

It is not possible to remove the bias and still get a shorter honest 90% interval than $\text{CI}_{\text{wide}} = \hat{\mu}_{\text{wide}} \pm 1.645 \hat{\tau}_{\text{wide}}/\sqrt{n}$. This follows from analysis similar to that in 4G. Thus, in a way, within the chosen large-sample framework, and provided we insist on guaranteed levels, we cannot carry out confidence and testing analysis better than with wide model methods, despite the fact that narrow estimators often have better precision than wide ones. A practical proposal is to use $\hat{\mu}_{\text{narr}}$ when theory and analysis suggest that it is more precise than $\hat{\mu}_{\text{wide}}$, but to supplement it with a confidence interval obtained through nonparametric or wide-model-parametric bootstrapping. The point is to obtain an honest 90% interval, for example, built around $\hat{\mu}_{\text{narr}}$.

Let us finally point out that narrow based intervals in some natural ways perform better than wide model ones, under mild misspecifications, since they are, indeed, narrower. Assume the loss incurred by using CI to cover μ is of the form

$$L\{(\theta, \gamma), \text{CI}\} = I\{\mu(\theta, \gamma) \notin \text{CI}\} + \sqrt{n}w \text{length}(\text{CI}),$$

where w is an appropriately chosen weight factor. The idea is to combine the two desiderata of confidence intervals into one measure; they should miss with low probability and have short length. The asymptotic risk functions for $\text{CI}_{\text{narr}} = \hat{\mu}_{\text{narr}} \pm z_0 \hat{\tau}_0/\sqrt{n}$ and $\text{CI}_{\text{wide}} = \hat{\mu}_{\text{wide}} \pm z_1 \hat{\tau}/\sqrt{n}$, under model P_n , become

$$\begin{aligned} \text{risk}_{\text{narr}} &= \Pr[|N\{b\delta/\tau_0, 1\}| \geq z_0] + 2wz_0\tau_0, \\ \text{risk}_{\text{wide}} &= \Pr[|N\{0, 1\}| \geq z_1] + 2wz_1(\tau_0^2 + b^2\kappa^2)^{1/2}. \end{aligned} \quad (4.4)$$

Again the best narrow method will be better than the best wide method, for moderate deviances δ from zero.

4I. Deviances in several directions. It could be worthwhile to study two types of model departure simultaneously, like both quadraticity and variance heterogeneity in regression. Generalising our results to such a situation is not difficult. Suppose $f(y, \theta_0, \gamma_0 + \delta/\sqrt{n})$ is the true model, where $\gamma = (\gamma_1, \gamma_2)'$ and $\delta = (\delta_1, \delta_2)'$ are two-dimensional. The natural criterion for when each narrow estimator is asymptotically more precise than its wide contender becomes $\delta\delta' \leq J^{22}$, where $J^{22} = (J_{22} - J_{12}'J_{11}^{-1}J_{12})^{-1}$ is 2×2 . This describes an ellipse around the null model. The border line, the crossing of which means coming into wide supremacy territory, is about $\text{Tr}(J_{22}J^{22})/2n$ away, as measured by Kullback–Leibler. The power of the 5% level $Z_n^2 = n(\hat{\gamma} - \gamma_0)'(\hat{J}^{22})^{-1}(\hat{\gamma} - \gamma_0)$ test is 33.4% at the border.

Most of the general results of Section 5 can similarly be extended to the situation with more than one type of deviance present. See remark (v) of 5H.

5. Classes of compromise estimators

We have so far concentrated on $\hat{\mu}_{\text{narr}}$ and $\hat{\mu}_{\text{wide}}$ to estimate $\mu = \mu_{\text{true}} = \mu(\theta_0, \gamma_0 + \delta/\sqrt{n})$. These rather cyclopic estimators can however be combined to form dimeric ones that perhaps work well both under the null model and the local alternative. This section considers and develops various more complex estimators with this aim. Some key words indicating the different types that will be discussed are pre-test or if-else estimators, mixture or weighted estimators, Bayes and empirical Bayes estimators, minimax estimators, the Bayesian epsilon estimator, and limited translation estimators.

Comparing all of these approaches may appear to be a formidable task, since the problem conceivably depends upon the particularities of the narrow model, the type and degree of deviance from it, and the specific parameter estimand under study. The comparison problem can however be drastically reduced, as we show in 5D below. Each of a large class of estimators for μ_{true} has a cousin which estimates a in a $N\{a, 1\}$ situation with one observation under squared error loss! The underlying one-one correspondence makes it possible to study the performance of general estimation approaches rather simply and rather generally, and this is indeed done in Section 6.

5A. If-else of pre-test estimators. ‘The responsibility of tolerance lies with those who have the wider vision.’ A procedure that is sometimes advocated in model choice problems and which perhaps is consistent with George Eliot’s view is as follows, in the present context: Test the hypothesis $\gamma = \gamma_0$ against the alternative $\gamma \neq \gamma_0$, say at the 10% level; if accepted, then use $\hat{\mu}_{\text{narr}}$, if rejected, then use $\hat{\mu}_{\text{wide}}$. Choosing the $Z_n^2 = n(\hat{\gamma} - \gamma_0)^2/\hat{\kappa}^2$ test also discussed in 4B, this suggestion amounts to

$$\hat{\mu}_{\text{pre}} = \hat{\mu}_{\text{narr}} I\{Z_n^2 \leq 1.645^2\} + \hat{\mu}_{\text{wide}} I\{Z_n^2 > 1.645^2\}, \quad 1.645^2 = \text{upper 10\% point of } \chi_1^2. \quad (5.1)$$

But this method sticks too rigidly to the narrow model. The theory of Section 3 suggests that one should rather use the much smaller value 1 as cut-off point, since $|\delta| \leq \kappa$ corresponds to $n(\gamma - \gamma_0)^2/\kappa^2 \leq 1$, and Z_n^2 estimates this ratio. Using 1 as cut-off is asymptotically the same as the Akaike model choice strategy, see 4D, and corresponds to a much more relaxed significance level, indeed to 31.7%, which in a way becomes the optimally chosen significance level. Observe that $\sqrt{n}(\hat{\mu}_{\text{pre}} - \mu_{\text{true}})$ tends to a mixture of two normals, as further commented upon below.

5B. Mixture estimators. Another natural idea is $\hat{\mu}_{\text{lin}} = (1 - c)\hat{\mu}_{\text{narr}} + c\hat{\mu}_{\text{wide}}$. To find the approximate distribution of this estimator it is necessary to go somewhat beyond the basic proposition of Section 2, in that the simultaneous limit distribution of the narrow and the wide estimators is needed. This can be found by studying the proof of the proposition, however. Utilising (2.4) and (2.5) it follows via some analysis that

$$\begin{aligned} \sqrt{n}(\hat{\mu}_{\text{narr}} - \mu_{\text{true}}) &\rightarrow_d b\delta + \left(\frac{\partial \mu}{\partial \theta}\right)' J_{11}^{-1} M, \\ \sqrt{n}(\hat{\mu}_{\text{wide}} - \mu_{\text{true}}) &\rightarrow_d \left(\frac{\partial \mu}{\partial \theta}\right)' J_{\text{wide}}^{-1} \begin{pmatrix} M \\ N \end{pmatrix}, \\ Z_n = \sqrt{n}(\hat{\gamma} - \gamma_0)/\hat{\kappa} &\rightarrow_d Z = (\delta + J^{21} M + J^{22} N)/\kappa, \end{aligned} \quad (5.2)$$

in which $(M, N) \sim N_{p+1}\{0, J_{\text{wide}}\}$. The convergence is simultaneous, and takes place under the P_n sequence of models (2.2). Note that $Z \sim N\{\delta/\kappa, 1\}$.

Now the limit distribution of $\hat{\mu}_{\text{lin}}$ can be obtained. The result is

$$\sqrt{n}(\hat{\mu}_{\text{lin}} - \mu_{\text{true}}) \rightarrow_d (1 - c)b\delta + (1 - c)\left(\frac{\partial \mu}{\partial \theta}\right)' J_{11}^{-1} M + c \left(\frac{\partial \mu}{\partial \gamma}\right)' J_{\text{wide}}^{-1} \begin{pmatrix} M \\ N \end{pmatrix}.$$

This is a normal distribution with mean value $(1 - c)b\delta$, and its variance can be shown to be $\tau_0^2 + c^2 b^2 \kappa^2$, in the notation of Section 3. The ideal value of c that minimises the asymptotic mean squared error for $\hat{\mu}_{\text{lin}}$ is $c_0 = \delta^2 / (\kappa^2 + \delta^2) = a^2 / (1 + a^2)$, featuring the key quantity $a = \delta/\kappa$. The accompanying minimum value is equal to $b^2 \kappa^2 a^2 / (1 + a^2) + \tau_0^2$. Note that this is always better than both the $b^2 \kappa^2 + \tau_0^2$ achieved by $\hat{\mu}_{\text{wide}}$ and the $b^2 \delta^2 + \tau_0^2$ achieved by $\hat{\mu}_{\text{narr}}$.

The problem is of course that c_0 is unknown since δ is. Using the empirical counterpart of $\delta = \sqrt{n}(\gamma - \gamma_0)$ invites $Z_n = \sqrt{n}(\hat{\gamma} - \gamma_0)/\hat{\kappa}$ to be inserted for δ/κ , i.e. Z_n^2 estimates a^2 , and one could try out the diophthalm

$$\hat{\mu}_{\text{eb}} = \frac{1}{1 + Z_n^2} \hat{\mu}_{\text{narr}} + \frac{Z_n^2}{1 + Z_n^2} \hat{\mu}_{\text{wide}}. \quad (5.3)$$

Note the Steinean overtones. The empirical Bayes connection that gives its subscript is noted in 5F below.

5C. Compromise estimators. Let us generalise. We shall be content to study estimators in the fairly large class of *compromise estimators*, which are bilingual and want the best from two worlds, and which we describe as follows. Its prime members are of the type

$$\mu^* = \{1 - c(Z_n)\} \hat{\mu}_{\text{narr}} + c(Z_n) \hat{\mu}_{\text{wide}}, \quad (5.4)$$

where $c(z)$ is almost everywhere continuous. Note that the previously considered estimators are of this form. The additional members that are included are those that can be closely approximated by (5.4) type ones by linearisation. More specifically, the limit distribution result (5.5) below is required to hold. It suffices for μ^* to be of the form $m(\hat{\mu}_{\text{narr}}, \hat{\mu}_{\text{wide}}, Z_n)$ for some smooth function $m(\mu_1, \mu_2, z)$ with the property that $m(\mu, \mu, z) \equiv \mu$. An example is the harmonic variety $\exp[\{1 - h(Z_n)\} \log \hat{\mu}_{\text{narr}} + h(Z_n) \log \hat{\mu}_{\text{wide}}]$ (which can be used in cases where μ is positive).

5D. Comparison of estimators: a drastic reduction. We wish to study the performance of all these estimators, and to compare pairs of them, w.r.t. the limiting mean squared error criterion.

THEOREM. *The compromise estimator (5.4) has limit distribution, under P_n of (2.2), given by*

$$\sqrt{n}(\mu^* - \mu_{\text{true}}) \rightarrow_d \Lambda = \{1 - c(Z)\} \{b\delta + \left(\frac{\partial \mu}{\partial \theta}\right)' J_{11}^{-1} M\} + c(Z) \left(\frac{\partial \mu}{\partial \gamma}\right)' J_{\text{wide}}^{-1} \begin{pmatrix} M \\ N \end{pmatrix}. \quad (5.5)$$

The mean squared error of the limit distribution can be written as

$$E\Lambda^2 = b^2 \kappa^2 E\{\delta/\kappa - c(Z)Z\}^2 + \tau_0^2 = b^2 \kappa^2 R(\delta/\kappa) + \tau_0^2, \quad (5.6)$$

in which

$$R(a) = E\{c(Z)Z - a\}^2 \text{ and } Z \sim N\{a, 1\}. \quad (5.7)$$

PROOF: (5.5) follows from (5.2) and the continuous mapping theorem of weak convergence. To characterise this limit variable Λ , study its distribution conditional on $Z = z$. Ordinary techniques from multivariate analysis, working from (5.2), lead to

$$\binom{M}{N} | \{Z = z\} \sim N_{p+1} \left\{ \begin{pmatrix} 0 \\ (\kappa z - \delta)/\kappa^2 \end{pmatrix}, \begin{pmatrix} J_{11} & J_{12} \\ J_{21} & J_{22} - 1/\kappa^2 \end{pmatrix} \right\}.$$

Several algebraic and multivariate details later one arrives at

$$\Lambda | \{Z = z\} \sim N\{b\delta - c(z)b\kappa z, \tau_0^2\}, \quad \text{where } Z \sim N\{\delta/\kappa, 1\}.$$

Expression (5.6) for the limiting mean squared error can now be worked out, studying first the z -conditional and then the unconditional mean value of Λ^2 . \square

This result contains those associated with (3.1) and (3.2) as well as (5.1) and the case of fixed c studied above. Observe that the unconditional distribution of Λ is non-normal unless $c(z)$ is constant in z . Note also that the unfamiliar type of limit distribution is not a peculiarity of the chosen local neighbourhood asymptotics, since Λ is typically non-normal even in the null model case.

A particular consequence of the theorem is that *it suffices to compare different versions of the function $R(a)$, as a function of $a = \delta/\kappa$, since $b\kappa$ and τ_0 remain unchanged for different estimators.* [We disregard the rather simple cases in which $b = 0$, see ‘case (i)’ of Section 3’s Result, under which all compromise estimators have $N\{0, \tau_0^2\}$ as limit distribution.] This constitutes an impressive reduction of the original comparison problem. Note that $R(a)$ is simply the risk function for the estimator $c(Z)Z$ for a in the one-observation $Z \sim N\{a, 1\}$ problem under squared error loss. There is a simple one-to-one correspondence from general compromise estimators to estimators $\hat{a}(Z)$ of a based on Z , via

$$\hat{a}(z) = c(z)z, \quad c(z) = \hat{a}(z)/z. \quad (5.8)$$

We stress the generality: A comparison between the four natural estimators $\hat{\mu}_{\text{narr}}$, $\hat{\mu}_{\text{wide}}$, $\hat{\mu}_{\text{pre}}$ of (5.1), and $\hat{\mu}_{\text{eb}}$ of (5.3), for example, can be carried out entirely in the realm of the classical $Z \sim N\{a, 1\}$ situation, by simply drawing the four $R(a)$ curves. See Section 6 for examples. And the conclusions from this comparison remain correct and relevant in *every* ‘moderate misspecification’ problem, cf. the wide span of problems that Examples A–G represent. Finally one is allowed to go the other way: *Your favourite estimator for a in the $N\{a, 1\}$ problem (where a may be rumoured to be in the vicinity of zero) can be transported to a useful estimator for any given estimand in any given moderate misspecification situation.*

In most cases it holds that $c(z) = c(-z)$, implying $R(a) = R(-a)$, making it necessary to study only non-negative a ’s. The parameter a measures the degree of misspecification from the narrow model. The important range is perhaps $[-4, 4]$, where $a = 0$ means correctness of the narrow model, $a = \pm 1$ are the turning points after which the wide estimator becomes better than the narrow one, and values beyond ± 3 could be thought of as clearly detectable departures from the narrow model, cf. power considerations (4.1), (4.2). [The 5% level test has power .851 and .979 at $a = 3$ and $a = 4$, whereas the 10% level test has .912 and .991 at the same points.]

These remarks also illustrate the importance of thinking about prior information related to the parameter a , for example its possible range. In Example B, studied in Sections 1 and 7 and in Hjort (1991a), a must be non-negative a priori, and in other cases it could be natural to restrict attention to the $[-4, 4]$ range, say, or to postulate a prior density for a . Such a prior could reflect serious prior beliefs, in the Bayesian fashion, or be used as a mathematical device to reach an estimator with minimum possible averaged mean squared error. Objectivists fretting at such ideas should note that the two classical solutions here, $\hat{\mu}_{\text{narr}}$ and $\hat{\mu}_{\text{wide}}$, correspond to full faith in the priors I_0 and 1, respectively, where I_0 is the degenerate distribution at zero and 1 is the flat non-informative prior for a . This is made clear in the course of the two following subsections, where the correspondence between moderate misspecification problems and the $N\{a, 1\}$ situation is explained also for Bayesian matters.

5E. Prior and posterior distributions for a . One is used to seeing that ‘the prior is washed out by the data’. Assume for example that a prior density $p_0(\theta, \gamma)$ is placed on (θ, γ) , with resulting Bayes estimators $(\hat{\theta}_B, \hat{\gamma}_B)$, expected values in the posterior density $p_0(\theta, \gamma|\text{data})$. Then these are typically asymptotically equivalent to the ML estimators, in the precise sense that $\sqrt{n}(\hat{\theta} - \hat{\theta}_B) \rightarrow_p 0$ and $\sqrt{n}(\hat{\gamma} - \hat{\gamma}_B) \rightarrow_p 0$, in the frequentist framework P_n . This is a fairly standard result under null model conditions, and the more delicate case of $\delta \neq 0$ can be treated using methods in Hjort (1986b).

This result uses a fixed prior for (θ, γ) , and is somewhat irrelevant in the present context of moderate misspecification. It appears more natural to operate with a fixed prior for $(\theta, \delta) = (\theta, \sqrt{n}(\gamma - \gamma_0))$, or, equivalently, a fixed prior $p(\theta, a)$ for $(\theta, a) = (\theta, \sqrt{n}(\gamma - \gamma_0)/\kappa)$. We think of the prior distribution for a as reflecting prior beliefs about the suitability of the narrow $f(y, \theta, \gamma_0)$ model, cf. the discussion above.

In this situation the prior information regarding θ will still be overwhelmed by the data, but not the part related to a . Information about a lies in $Z_n = \sqrt{n}(\hat{\gamma} - \gamma_0)/\hat{\kappa}$, which is not consistent, but has a limiting variable $Z \sim N\{a, 1\}$. Intuitively, therefore, the posterior density $p(a|\text{data})$ should for large n simply be close to $p(a|z)$ in the situation where Z is $N\{a, 1\}$ and a has prior proportional to $p(\theta_0, a)$. To prove it, let us study $p(a|Y_1, \dots, Y_n)$ when n grows, under the P_n model, where $f(y) = f(y, \theta_0, \gamma_0 + \delta_0/\sqrt{n})$ for some fixed values of θ_0, δ_0 . Let $L_n(\theta, \gamma) = \prod_{i=1}^n f(Y_i, \theta, \gamma)$ be the n 'th likelihood. By judicious second order Taylor expansion analysis it can be established that

$$H_n(s, t) = \frac{L_n(\hat{\theta} + s/\sqrt{n}, \hat{\gamma} + t/\sqrt{n})}{L_n(\hat{\theta}, \hat{\gamma})} \rightarrow_d H(s, t) = \exp\left\{-\frac{1}{2} \begin{pmatrix} s \\ t \end{pmatrix}' J_{\text{wide}} \begin{pmatrix} s \\ t \end{pmatrix}\right\}$$

under P_n of (2.2). The convergence takes place in each Скороход space $D[-A, A]^{p+1}$. Let now $g(\theta, a)$ be any bounded function. Then one may deduce

$$\begin{aligned} E\{g(\theta, a)|\text{data}\} &= \frac{\int \int g(\theta, \sqrt{n}(\gamma - \gamma_0)/\kappa) L_n(\theta, \gamma) p(\theta, \sqrt{n}(\gamma - \gamma_0)/\kappa) \sqrt{n}/\kappa d\theta d\gamma}{\int \int L_n(\theta, \gamma) p(\theta, \sqrt{n}(\gamma - \gamma_0)/\kappa) \sqrt{n}/\kappa d\theta d\gamma} \\ &= \frac{\int \int g(\hat{\theta} + s/\sqrt{n}, Z'_n + t/\kappa) H_n(s, t) p(\hat{\theta} + s/\sqrt{n}, Z'_n + t/\kappa) ds dt}{\int \int H_n(s, t) p(\hat{\theta} + s/\sqrt{n}, Z'_n + t/\kappa) ds dt} \end{aligned}$$

$$\begin{aligned}
& \rightarrow_d \frac{\int \int g(\theta_0, Z + t/\kappa) H(s, t) p(\theta_0, Z + t/\kappa) ds dt}{\int \int H(s, t) p(\theta_0, Z + t/\kappa) ds dt} \\
& = \frac{\int g(\theta_0, Z + t/\kappa) \exp(-\frac{1}{2}t^2/\kappa^2) \pi(Z + t/\kappa) dt}{\int \exp(-\frac{1}{2}t^2/\kappa^2) \pi(Z + t/\kappa) dt} \\
& = \frac{\int g(\theta_0, a) \exp\{-\frac{1}{2}(Z - a)^2\} \pi(a) da}{\int \exp\{-\frac{1}{2}(Z - a)^2\} \pi(a) da},
\end{aligned}$$

in which $\pi(a) = \text{const. } p(\theta_0, \delta)$ is the prior for a given the information $\theta = \theta_0$, and $Z'_n = (\hat{\kappa}/\kappa)Z_n$ was used for notational simplicity. The necessary mathematical details have to do with (i) securing convergence inside $[-A, A]^{p+1}$, utilising the proposition of Section 2, along with (5.2); (ii) cleverly carrying out a certain inner p -dimensional normal integration; and (iii) bounding integrands outside $[-A, A]^{p+1}$ for large A . The arguments that are needed resemble those explained in Hjort (1986b) (to reach a different conclusion, in a different problem), and are left out here.

By considering $g = g(a)$ above it is clear that

$$\pi_n(a|\text{data}) \rightarrow_d \pi(a|Z) = \frac{\phi(Z - a)\pi(a)}{\int \phi(Z - a)\pi(a) da}, \quad (5.9)$$

under P_n , where $Z \sim N\{a, 1\}$ is as in (5.2). This is what was predicted above. If a has some prior distribution $d\pi(a)$ that perhaps does not have a density, then the arguments can be repeated to give $d\pi_n(a|\text{data}) \rightarrow_d \text{const. } \phi(Z - a) d\pi(a)$ instead.

5F. Bayes and empirical Bayes estimators. We should distinguish between kosher Bayes and approximate Bayes estimators. A prior density $p(\theta, a)$ for (θ, a) leads to the exact Bayes solution $\hat{a}_n = E\{a|Y_1, \dots, Y_n\}$. This is usually a very complicated expression, and in view of (5.9) it is tempting to work directly in the limit distribution and use $\hat{a}(Z_n)$ instead, where

$$\hat{a}(z) = E\{a|Z = z\} = \frac{\int a\phi(z - a)\pi(a) da}{\int \phi(z - a)\pi(a) da} = z + \frac{\partial}{\partial z} \log \int \phi(z - a)\pi(a) da. \quad (5.10)$$

But the arguments of 5E can be used to reach

$$\hat{a}_n = Z'_n + \frac{\partial}{\partial z} \log \int \phi(Z'_n - a)\pi(a) da + O_p(1/\sqrt{n}),$$

where again $Z'_n = (\hat{\kappa}/\kappa)Z_n$. This proves $\hat{a}_n - \hat{a}(Z_n) \rightarrow_p 0$, under P_n , allowing us to use $\hat{a}(Z_n)$ instead when we devise and study Bayes solutions in our large-sample framework. In particular we do not have to bother with the part of the prior information that has to do with θ .

Some specific Bayesian and empirical Bayesian constructions follow.

(i) Suppose a is $N\{0, \sigma^2\}$ (where the size of the spread parameter σ matters more than the normality). Then $Z \sim N\{0, \sigma^2 + 1\}$, and $\hat{a}(z) = \{\sigma^2/(\sigma^2 + 1)\}z$, with Bayes risk

$\sigma^2/(\sigma^2 + 1)$. If $Ea^2 = \sigma^2$ is unknown, a simple guess is Z_n^2 , since Z_n estimates a . This brings forward the empirical Bayes estimate $\hat{a}_{\text{eb}}(Z_n) = \{Z_n^2/(Z_n^2 + 1)\} Z_n$ for a . But this corresponds to $\hat{\mu}_{\text{eb}}$ of (5.3), explaining its empirical Bayes interpretation.

One may also consider other estimators for $q = \sigma^2/(\sigma^2 + 1)$ here. Each such $\hat{q} = \hat{q}(Z)$ leads to an $a^* = \hat{a}(Z, \hat{q})$, and in its turn to a new estimator μ^* for μ_{true} via (5.4) and (5.8). The fact that $EZ^2 = \sigma^2 + 1$ suggests $\tilde{q} = (Z^2 - 1)_+/Z^2$, which in fact is the ML solution, or similar versions. Another proposal is to put a vague hyper prior on σ , or directly on the ratio q . The Bayes solution becomes $E\{a|Z = z\} = E_z E_z\{a|\sigma\} = E_z(qz) = \hat{q}(z)z$, in which $\hat{q}(z) = \int_0^1 qp(q|z) dq$ is the posterior density of q for given $Z = z$. The usual choice for a non-informative prior for a scale parameter like σ is to have $\log \sigma$ uniform. This leads to $p(q) = \text{const.} \{q(1 - q)\}^{-1}$ on $[\varepsilon, 1 - \varepsilon]$, say, for q , with a corresponding explicit $\hat{q}(z)$. In fact it turns out that

$$\hat{q}(z) = \frac{\int_{\varepsilon}^1 (1 - q)^{-1/2} \exp\{-\frac{1}{2}(1 - q)z^2\} dq}{\int_{\varepsilon}^1 q^{-1}(1 - q)^{-1/2} \exp\{-\frac{1}{2}(1 - q)z^2\} dq} \quad (5.11)$$

is substantially better than $\tilde{q} = (z^2 - l)_+/(z^2 + 1 - l)$, where $0 \leq l \leq 1$, for a wide right interval $(q_0, 1)$ of q values. Heroic numerical integrations have demonstrated this, via computations and comparisons of $E_q|q^* - q|$ for the various estimators. The \hat{q} above, with $\varepsilon = 0.05$, is for example much better than the \tilde{q} ones, for q in $(0.20, 1)$.

(ii) Suppose a comes from $\pi_0(a)$ with probability p_0 and from $\pi_1(a)$ with probability p_1 . Then $\hat{a}(z)$ can be shown to be of the mixture form $w_0(z)\hat{a}_0(z) + w_1(z)\hat{a}_1(z)$, where $\hat{a}_j(z)$ is the Bayes estimator under theory $p_j(a)$, and $w_j(z) = p_j h_j(z)/\{p_0 h_0(z) + p_1 h_1(z)\}$, and $h_j(z) = \int \phi(z - a)\pi_j(a) da$. An interesting special case is the prior distribution $a \sim (1 - \varepsilon)I_0 + \varepsilon N\{0, \sigma^2\}$, where I_0 denotes the degenerate distribution at zero. This is a ‘Bayesian epsilon’ approach, where the statistician is rather convinced of the narrow model’s correctness but allows the data to express a different opinion with probability ε . In this case

$$\hat{a}(z) = \frac{\varepsilon}{\varepsilon + (1 - \varepsilon)B(z)} \frac{\sigma^2}{\sigma^2 + 1} z, \quad B(z) = \frac{h_0(z)}{h_1(z)} = \sqrt{\sigma^2 + 1} \exp\left\{-\frac{1}{2} \frac{\sigma^2}{\sigma^2 + 1} z^2\right\}. \quad (5.12)$$

Again σ^2 has to be specified or estimated. One possibility is $\hat{\sigma}^2 = Z^2/\varepsilon$, since $Ea^2 = \varepsilon\sigma^2$ and Z estimates a ; other versions can be constructed as in (i) above.

(iii) If it is assumed that $|a| \leq m$ a priori then the Bayes solution (5.10) with a uniform prior on $[-m, m]$ should give an estimator with good risk properties on this interval.

5G. Minimax type estimators. The remarks about the a parameter in 5D suggest that its range could usefully be taken to be bounded, a priori, in some situations. If a is postulated to be in $[-m, m]$, for example, then estimators a^* exist that are uniformly better than z , which means, by our basic correspondence theorem, that estimators μ^* exist that are uniformly better than $\hat{\mu}_{\text{wide}}$. If in particular a_m^* is a minimax estimator, with maximum risk $r_m < 1$ for $R(a)$ in $[-m, m]$, then μ^* of (5.4), defined via (5.8), has a minimax property: It minimises the limit distribution version of

$$\max_{|\gamma - \gamma_0| \leq m\kappa/\sqrt{n}} nE_{\theta_0, \gamma}\{\hat{\mu} - \mu(\theta_0, \gamma)\}^2$$

over all estimators $\hat{\mu}$, and achieves $\max_{|\delta| \leq m\kappa} E\Lambda^2 = b^2\kappa^2 r_m + \tau_0^2 < \tau^2$.

How do such minimax estimators look like? It is known that $a_m^*(z)$ is the proper Bayes solution w.r.t. a prior distribution concentrated in a finite number of points, see e.g. Lehmann (1983, Chapter 4.3). This least favourable prior has been found for small values of m , at least for $m \leq 1.5$, and Bickel (1981) gives approximate results for m large. We mention that $a^* = m \tanh(mz)$, the Bayes solution under a symmetric two-point prior in $\pm m$, is minimax, provided $m \leq 1.05$. This is relevant here since $[-1, 1]$ is the range for a where narrow estimation is better than wide estimation. Bickel shows that the distribution with density $\pi_m(a) = \cos^2(\frac{1}{2}\pi a/m)/m$ for $|a| \leq m$ is approximately least favourable, for large m . This suggests trying out

$$\hat{a}_{\text{bic}}(z) = \frac{\int_{-m}^m a \phi(z-a) \cos^2(\frac{1}{2}\pi a/m) da}{\int_{-m}^m \phi(z-a) \cos^2(\frac{1}{2}\pi a/m) da}, \quad (5.13)$$

It is *not* approximately minimax on $[-m, m]$, but it is uniformly better than $\hat{a}_{\text{wide}} = z$ in a certain interval around 0. A simpler possibility is to use the ML solution when $|a| \leq m$ a priori, that is,

$$\hat{a}_{\text{res}}(z) = -m \text{ when } z \leq -m, \quad z \text{ on } [-m, m], \quad m \text{ when } z \geq m. \quad (5.14)$$

This is not quite as good as using the proper minimax solution on $[-m, m]$, however.

Finally we should include estimators of the Efron–Morris variety, see Efron and Morris (1971) and Lehmann (1986, Chapter 4.2). These aim at minimising Bayes risk, under normal priors, subject to having maximum risk less than some prescribed level. A particular case of these is pertinent here, namely the ‘limited translation estimator’

$$\hat{a}_{\text{em}}(z) = z + m \text{ when } z \leq -m, \quad 0 \text{ on } [-m, m], \quad z - m \text{ when } z \geq m. \quad (5.15)$$

These come close to minimising maximum risk subject to doing well at $a = 0$, see Bickel (1983, 1984) and Berger (1982). They are not smooth enough to be admissible. An alternative estimator which can be proposed is

$$\hat{a}_{\text{atan}}(z) = z - m(2/\pi) \arctan z. \quad (5.16)$$

It is motivated from Bickel’s study of \hat{a}_{em} and its connection to bounded influence functions in robust estimation of location, and is scaled so that it has the same maximum risk $1 + m^2$ as (5.15) (see below).

5H. Some concluding comments.

(i) Observe the generality under which the comparisons of Sections 5 and 6 are made. They are valid and relevant for all of Examples A–G (with appropriate modifications for case B, see Hjort (1991a)) and for all parameter estimands, via (5.4)–(5.7).

(ii) When applied to a particular estimand in a particular model these comparisons should perhaps also include the nonparametric contender. In Example A, for example, one could compare the wide and the narrow parametric methods to the sample median.

(iii) We have been motivated by approximate mean squared error $E_{P_n}(\mu^* - \mu_{\text{true}})^2$ when comparing estimator performance. We haven't quite worked with the limit of n times the mean squared error, but rather with $E\Lambda^2$ in (5.6), using the limit distribution. This is both easier and more meaningful. This is a minor technical point, however; usually the two agree. See Lehmann's (1983) Lemma 5.1.2, for example.

(iv) One might wish to study L_1 error $\sqrt{n}E_{P_n}|\mu^* - \mu_{\text{true}}|$ and its limit distribution version $E|\Lambda|$ instead. There is a parallel result to (5.6) and (5.7) for this problem. Let $L(x)$ be the function $E|x + N\{0, 1\}| = x + 2\phi(x) - 2x\{1 - \Phi(x)\}$. Then μ^* of (5.4) has

$$E|\Lambda| = \tau_0 \int L((b\kappa/\tau_0)(\delta/\kappa - c(z)z)) \phi(z - a) dz = \tau_0 E_a L(\rho(c(Z)Z - a)), \quad (5.17)$$

letting $a = \delta/\kappa$ again and $\rho = |b|\kappa/\tau_0$. There is once more a correspondence between compromise estimators μ^* of μ and estimators $\hat{a}(z) = c(z)z$ of a , but the L_1 loss function $|\mu^* - \mu|$ for μ is transformed to loss function $L(\rho(\hat{a} - a))$ for a . And there is still a 'tolerance radius' around the narrow model inside of which misspecification is favourable, but one does not get the clear-cut $|\delta| \leq \kappa$ answer. The narrow and the wide procedures have respectively $E|b\delta + \tau_0 N|$ and $(b^2\kappa^2 + \tau_0^2)^{1/2} E|N|$ as limiting risks, where $N \sim N\{0, 1\}$. The tolerance radius becomes in fact $|\delta| \leq a_0\kappa = a_0(\rho)\kappa$, or $|a| \leq a_0(\rho)$, where $|a| \leq a_0$ corresponds to $E|\rho a + N| \leq (1 + \rho^2)^{1/2} E|N|$. Computations show that $a_0(\rho)$ starts at 1.00 for $\rho = 0$ and slouches towards $\sqrt{2/\pi} = 0.7979$ as ρ grows. The L_1 criterion for estimation of μ accordingly tolerates slightly less misspecification than the L_2 criterion.

(v) One might generalise most of this section's results to deviances from the basic model in more than one direction. One can envisage useful generalisations of (5.4), for example, where the final estimator gives weights to the narrow model and to several wider alternative models, with weights determined by the data. Statistics tradition does perhaps dictate this point of view, with a classic null model and several possible departures from it, but the problem can also be turned inside out, starting with a wide a priori model for the data and then smoothing in several directions downwards to narrower models of interest. The empirical Bayes ideas above should be of value if these kind of questions are to be pursued.

6. Grand comparison

Each estimator of μ_{true} has a cousin that estimates a in the $Z \sim N\{a, 1\}$ situation, and vice versa, by (5.8). Furthermore, the performance of one of them determines and is determined by the performance of the other one, by the key correspondence (5.6)–(5.7). It is refreshing to judge a μ -estimator by examining its a -estimator cousin. Here is a partial list of interesting estimators for μ_{true} , following the various suggestions of Section 5, along with brief descriptions of their performance.

(i) The narrow estimator $\hat{\mu}_{\text{narr}}$ has $c(z) \equiv 0$ and $\hat{a}(z) \equiv 0$. This particular estimator of a is fully confident that a is close to zero, and has risk $R_{\text{narr}}(a) = a^2$.

(ii) The wide estimator $\hat{\mu}_{\text{wide}}$ on the other hand has $c(z) \equiv 1$ and $\hat{a}(z) = z$. This conservative estimator has constant risk $R_{\text{wide}}(a) = 1$, and is the unique admissible minimax estimator for a when the parameter range is unrestricted. Note anew that the narrow is better than the wide when $|a| \leq 1$.

(iii) The if-else estimator (5.1), with m^2 instead of 1.645^2 as cut-off point, has $c(z) = I\{|z| \geq m\}$, and corresponds to the a -estimator $\hat{a}(z) = zI\{|z| \geq m\}$. A determined mind finds

$$\begin{aligned} R_{\text{pre}}(a) &= \int_{|z| \geq m} (z - a)^2 \phi(z - a) dz + \int_{|z| \leq m} (0 - a)^2 \phi(z - a) dz \\ &= 1 + (a^2 - 1)\{\Phi(m + a) + \Phi(m - a) - 1\} \\ &\quad + (m + a)\phi(m + a) + (m - a)\phi(m - a). \end{aligned}$$

The if-else with cut-off $m = 1$, which corresponds to the Akaike strategy, see 4D, seems overall to be preferable to the one with $m = 1.645$, corresponding to the 10% level test. The latter is better in the vicinity of the narrow model, for $|a| \leq 0.83$, but then becomes markedly worse than the former. The pre-test estimators are not smooth enough to be Bayes or extended Bayes, see (5.10). In particular such methods are not admissible, i.e. they can be improved upon uniformly in a ! Note that $m = 0$ and $m = \infty$ give back the wide and the narrow methods, respectively. These extreme cases *are* admissible, however.

(iv) The linear combination estimator $\hat{\mu}_{\text{lin}}$ discussed in 5B has $c(z) = c$ and $\hat{a}(z) = cz$. Its risk is $R_{\text{lin}}(a) = c^2 + (1 - c)^2 a^2$, which is unbounded when $|a|$ grows. These are proper Bayes solutions for $0 \leq c < 1$ and admissible for $0 \leq c \leq 1$.

(v) The natural $\hat{\mu}_{\text{eb}}$ of (5.3) has $c(z) = z^2/(1 + z^2)$, and the correspondence to the empirical Bayes estimator $\hat{a}_{\text{eb}}(z) = \{z^2/(1 + z^2)\}z$ has been noted in 5F. One must compute

$$R_{\text{eb}}(a) = E_a \left\{ \frac{Z^3}{1 + Z^2} - a \right\}^2 = \int \left\{ \frac{z^3}{1 + z^2} - a \right\}^2 \phi(z - a) dz$$

by numerical integration. I can prove that \hat{a}_{eb} is admissible. This translates into an admissibility property for $\hat{\mu}_{\text{eb}}$. I have also studied the similarly inspired $\hat{a}(z) = \hat{q}(z)z$, with $\hat{q}(z)$ as in (5.11) instead of $z^2/(z^2 + 1)$. These perform similarly. The risk for $\hat{q}(z)z$ starts at 0.37 for $a = 0$, smaller than 0.46 for $\hat{a}_{\text{eb}}(z)$, and stays better for $|a| \leq 1.83$, after which $\hat{a}_{\text{eb}}(z)$ takes over. The maximum risk 1.476 for $\hat{q}(z)z$ is higher than 1.252 for $\hat{a}_{\text{eb}}(z)$. The risk for $\hat{a}_{\text{eb}}(z)$ is less than the crucial value 1 for $|a| \leq 1.40$. Overall one would argue that $\hat{a}_{\text{eb}}(z)$ is better than $\hat{q}(z)z$; see also the figure below.

(vi) The restricted ML estimator (5.14) can be shown to have risk function

$$\begin{aligned} R_{\text{res}}(a) &= \Phi(m - a) + \Phi(m + a) - 1 - (m - a)\phi(m - a) - (m + a)\phi(m + a) \\ &\quad + (m - a)^2\{1 - \Phi(m - a)\} + (m + a)^2\{1 - \Phi(m + a)\}. \end{aligned}$$

This estimator is not smooth enough to be Bayes or extended Bayes, and is like the if-else estimator not admissible. Its risk is satisfactory on $[-m, m]$ but ends up growing as a^2 outside it.

(vii) The Efron-Morris estimator (5.15) has risk function

$$R_{\text{em}}(a) = 1 + m^2 + (a^2 - m^2 - 1)\{\Phi(m + a) + \Phi(m - a) - 1\} - (m - a)\phi(m + a) - (m + a)\phi(m - a).$$

These increase with $|a|$, and rather rapidly, from a small value at zero towards maximum risk $1 + m^2$. The arctan-estimator (5.16) has also risk that increases in $|a|$ from $R_{\text{atan}}(0)$ to $1 + m^2$. It has higher risk than (5.15) has at $a = 0$, but the risk climbs much more slowly towards $1 + m^2$; see also the figure below. In particular an arctan-estimator can be better than an Efron-Morris estimator on $[-5, 5]$, say.

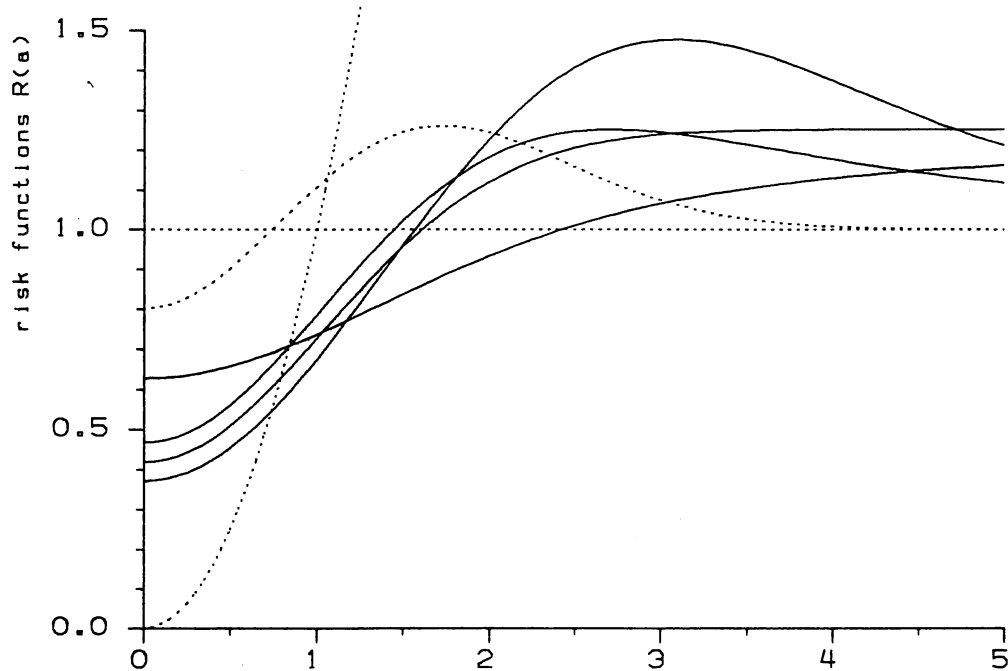


FIGURE. Risk functions $R(a)$ are shown for seven procedures, corresponding to seven choices of $c(Z_n)$ in (5.4). Risks for the wide and the narrow methods start at 1.00 and 0.00, and are shown with dotted lines, as is the risk for Akaike, which starts at 0.80. The empirical Bayes methods $\hat{a}_{eb}(Z)$ and $\hat{q}(Z)Z$ start at 0.47 and 0.37. Finally the Efron–Morris and arctan estimators, both scaled to have the same maximum risk 1.252 as has $\hat{a}_{eb}(Z)$, start at 0.42 and 0.63.

(viii) The Bickel-inspired estimator (5.13) has acceptable risk below 1 in an interval around 0, but the risk explodes when $|a|$ grows. The same goes for the Bayes solution with a uniform prior on $[-m, m]$. Its risk is 1 at 0 and at m and below 1 in between, but quickly explodes when $|a|$ grows outside the interval. Note that for m large this solution becomes simply the wide solution.

(ix) The ‘epsilon Bayes’ methods described in (5.12) and the remarks following it have small risk for $|a|$ less than about 1, but then become markedly worse than both $\hat{a}_{eb}(Z)$, $\hat{q}(Z)Z$, and pre-test estimators. The empirical epsilon Bayes method with $\hat{\sigma}^2 = Z^2/\varepsilon$ is not as good as the simple specified one with σ put equal to 3, for example.

Let us compare $\hat{\mu}_{narr}$, $\hat{\mu}_{wide}$, the if-else $\hat{\mu}_{pre}$ with $m = 1$, and the mixture estimator $\hat{\mu}_{eb}$ of (5.3). The narrow estimator wins if $|a| \leq 0.84$; the mixture estimator wins when $|a|$ is between 0.84 and 1.45, and finally the safest and wide estimator wins if $|a|$ exceeds 1.45. While $\hat{\mu}_{narr}$ can misbehave significantly when $|\delta| \geq 2.50\kappa$, say, $\hat{\mu}_{eb}$ always behaves wisely, also in the $|\delta| > \kappa$ case, and does not ever lose much to $\hat{\mu}_{wide}$. Its worst risk value is 1.252, at $|a| = 2.70$, and when the narrow model is very wrong ($|\delta|$ is large) $\hat{\mu}_{eb}$ becomes equivalent to $\hat{\mu}_{wide}$. In no region does $\hat{\mu}_{pre}$ win, but its risk function lies between the wide method’s 1 and the mixture method’s risk function, for $|a| > 2.17$; see the figure.

Based on these observations five of the more interesting estimators are singled out for display, in addition to the extreme basis choices ‘narrow’ and ‘wide’. The five are the empirical Bayes versions $\hat{a}_{eb}(Z)$ and $\hat{q}(Z)Z$; the Akaike strategy $\hat{a}_{pre}(Z)$, see 4D and 5A; the Efron–Morris (5.15) with $m = 0.502$ chosen so as to get the same maximum risk

1.252 as $\hat{\mu}_{\text{eb}}(Z)$; and the smoother arctan-estimator (5.16) with the same m (and the same objective). The Akaike strategy is about as good as the pre-test method can be, but is not as good as the others. It is included since versions of it are in frequent use.

All in all the best choices seem to be the simple empirical Bayes, the Efron–Morris, and the arctan. There are several other methods among those discussed that would make a good show on $[-5, 5]$, say, but with risks that explode for growing $|a|$. The $\hat{\mu}_{\text{eb}}$ of (5.3) in particular is a practical and satisfactory solution. There is no artificial cut-off, its weight in favour of the wide model is smoothly increasing from 0 to 1 with the test indicator Z_n , it behaves considerably better than the wide estimator in a reasonable neighbourhood of the narrow model, and its maximum risk is only $(1.119 b\kappa)^2 + \tau_0^2$, compared to $(b\kappa)^2 + \tau_0^2$ for the conservative wide method. The Efron–Morris and the arctan estimators have similar performances but require selection of a parameter, related to the trade-off between behaving well around zero and having a small maximum risk.

To illustrate more concretely what these suggestions amount to, consider logistic regression as in Example F. If deviation from $\alpha + \beta x$ in direction of quadraticity is suspected, use

$$p^*(x) = \frac{1}{1 + Z_n^2} \frac{\exp(\hat{\alpha}_{\text{narr}} + \hat{\beta}_{\text{narr}} x)}{1 + \exp(\hat{\alpha}_{\text{narr}} + \hat{\beta}_{\text{narr}} x)} + \frac{Z_n^2}{1 + Z_n^2} \frac{\exp(\hat{\alpha} + \hat{\beta} x + \hat{\gamma} x^2)}{1 + \exp(\hat{\alpha} + \hat{\beta} x + \hat{\gamma} x^2)},$$

where $Z_n = \sqrt{n}\hat{\gamma}/\hat{\kappa}$. Or replace the weights with $1 - \hat{a}(Z_n)/Z_n$ and $\hat{a}(Z_n)/Z_n$, with $\hat{a}(Z_n)$ equal to the limited translation estimator (5.15) or the arctan-estimator (5.16).

The facts above are meant to summarise the main features of the various estimator performances, based on a thorough investigation and several days of conscientious staring at hundreds of risk functions. Computer programs and risk tables for each of the estimator classes discussed above are collected in Hjort (1991b), which is available upon courteous request.

7. Examples

We now provide answers to the questions asked in Examples A–G of the introduction!

EXAMPLE A. In the general two-parameter Weibull model, parameterised as in (1.1), the score function becomes

$$\left(\begin{array}{c} \frac{\gamma}{\theta} \{1 - (\theta y)^\gamma\} \\ \frac{1}{\gamma} \{1 + \log(\theta y)^\gamma - (\theta y)^\gamma \log(\theta y)^\gamma\} \end{array} \right),$$

and clever computations involving the gamma function and its derivatives reveal the information matrix and its inverse to be

$$J_{\text{gen}} = \begin{pmatrix} \gamma^2/\theta^2 & (1-k)/\theta \\ (1-k)/\theta & c^2/\gamma^2 \end{pmatrix}, \quad J_{\text{gen}}^{-1} = \frac{1}{\pi^2/6} \begin{pmatrix} c^2\theta^2/\gamma^2 & -(1-k)\theta \\ -(1-k)\theta & \gamma^2 \end{pmatrix},$$

in which $k = 0.577\dots$ is the Euler–Mascheroni constant and $c^2 = \pi^2/6 + (1-k)^2$. The null model corresponds to $\gamma = \gamma_0 = 1$. The κ^2 parameter is $6/\pi^2$, and we have reached the following conclusion: For $|\gamma - 1| \leq \sqrt{6/\pi^2}/\sqrt{n} = 0.779/\sqrt{n}$, estimation with $\mu(1/\bar{Y}, 1)$ based on simple and narrow-minded exponentiality performs better than high-brow $\mu(\hat{\theta}, \hat{\gamma})$; and this is true regardless of the parameter μ to be estimated.

In the language of 4F Weibull deviance from exponentiality has tolerance limit $d = J_{22}J^{22} = 1 + (1-k)^2/(\pi^2/6) = 1.109$, and ρ^2 of (4.3) is $(1-k)^2/\{(1-k)^2 + \pi^2/6\} = 0.098$. It is instructive to compare these with corresponding values for gamma distribution deviance from exponentiality. If $f(y) = \{\theta^\gamma/\Gamma(\gamma)\} y^{\gamma-1} e^{-\theta y}$ is the gamma density, for which $\gamma_0 = 1$ gives back exponentiality, then $\kappa^2 = 1/(\pi^2/6 - 1)$; estimation using $\mu(1/\bar{Y}, 1)$ is more precise than $\mu(\hat{\theta}, \hat{\gamma})$ provided $|\gamma - 1| \leq 1.245/\sqrt{n}$; d is 2.551; and $\rho^2 = 6/\pi^2 = 0.608$. This suggests that moderate gamma-ness is less critical than moderate Weibull-ness for standard methods based on exponentiality.

EXAMPLE B. The wide model has parameters ξ, σ, m . Let us reparameterise to $\gamma = 1/m$, so that the density becomes

$$f(y, \xi, \sigma, \gamma) = \frac{c(\gamma)}{\sigma} \left\{ 1 + \gamma \left(\frac{y - \xi}{\sigma} \right)^2 \right\}^{-\{1/2 + 1/(2\gamma)\}}, \quad c(\gamma) = \frac{\sqrt{\gamma} \Gamma(1/2 + 1/(2\gamma))}{\sqrt{\pi} \Gamma(1/(2\gamma))}.$$

Estimation of the model parameters must now be studied when γ is small and nonnegative. This actually calls for special treatment since the null point $\gamma = 0$ is not an inner point, and $\hat{\gamma} = 0$, or $\hat{m} = \infty$, happens with positive probability. Such a treatment is given in Hjort (1991a), and shows that if $\gamma \leq 0.686/\sqrt{n}$, i.e. if the degrees of freedom $m \geq 1.458\sqrt{n}$, then t -ness doesn't matter, and *any* parameter $\mu = \mu(f) = \mu(\xi, \sigma, m)$ is better estimated in the ordinary, simple, normality based way. A similar result is also proven there for regression models.

EXAMPLE C. We generalise slightly and write the wide model as $Y_i \sim N\{\beta'x_i + \gamma c(x_i), \sigma^2\}$, where β and x_i are p -dimensional vectors, and $c(x)$ is some given scalar function. By computing log-derivatives and evaluating covariances one reaches

$$J_{n, \text{wide}} = \frac{1}{\sigma^2} \begin{pmatrix} 2 & 0 & 0 \\ 0 & n^{-1} \sum_{i=1}^n x_i x_i' & n^{-1} \sum_{i=1}^n x_i c(x_i) \\ 0 & n^{-1} \sum_{i=1}^n x_i' c(x_i) & n^{-1} \sum_{i=1}^n c(x_i)^2 \end{pmatrix},$$

from the definition in (3.4). It follows that

$$\kappa^2 = \sigma^2 \times \text{lower right element of } \begin{pmatrix} n^{-1} \sum_{i=1}^n x_i x_i' & n^{-1} \sum_{i=1}^n x_i c(x_i) \\ n^{-1} \sum_{i=1}^n x_i' c(x_i) & n^{-1} \sum_{i=1}^n c(x_i)^2 \end{pmatrix}^{-1}.$$

Assume, for a concrete example, that x_i is one-dimensional and uniformly distributed over $[0, b]$, say $x_i = b \frac{i}{n+1}$, and that the wide model has $\alpha + \beta(x_i - \bar{x}) + \gamma(x_i - \bar{x})^2$. Then $\kappa \doteq \sqrt{80} \sigma / b^2$. Consequently, dropping the quadratic term does not matter, and is actually advantageous, for every estimator, provided $|\gamma| \leq 8.94 \sigma / b \sqrt{n}$. In many situations with moderate n this will indicate that it is best to keep the narrow model and avoid quadratic analysis.

EXAMPLE D. Again we are mildly general and write $Y_i \sim N\{\beta'x_i, \sigma^2(1 + \gamma c(x_i))\}$ for the $p + 2$ parameter variance heterogeneous model. It is not easy to put up simple expressions for the general information matrix, in the presence of γ , but once more we are

permitted to compute $J_{n,\text{wide}}$ and J_{wide} of (3.4) under the null model, that is, when $\gamma = 0$. Some calculations give

$$J_{n,\text{wide}} = \begin{pmatrix} 2/\sigma^2 & 0 & \sigma^{-1}n^{-1} \sum_{i=1}^n c(x_i) \\ 0 & \sigma^{-2}n^{-1} \sum_{i=1}^n x_i x_i' & 0 \\ \sigma^{-1}n^{-1} \sum_{i=1}^n c(x_i) & 0 & (2n)^{-1} \sum_{i=1}^n c(x_i)^2 \end{pmatrix}.$$

Matters simplify and $1/\kappa^2$ is found to be $(2n)^{-1} \sum_{i=1}^n \{c(x_i) - \bar{c}\}^2$. If once again x_i 's are distributed evenly on $[0, b]$, and $c(x_i) = x_i$, then $\kappa \doteq \sqrt{24}/b$, and the criterion becomes $|\gamma| \leq 4.90/b\sqrt{n}$. In particular this shows that the sophisticated variance heterogeneous approach, which uses the weighted least squares estimator

$$\hat{\beta}_{\text{soph}} = \sum_{i=1}^n \frac{x_i y_i}{1 + \hat{\gamma} c(x_i)} / \sum_{i=1}^n \frac{x_i^2}{1 + \hat{\gamma} c(x_i)},$$

is inferior to the simpler solution, unless $|\gamma|$ is quite large. It is of course the sampling variability present in the weights, via the ML estimator $\hat{\gamma}$, that makes $\hat{\beta}_{\text{soph}}$ inferior to ordinary $\hat{\beta}_{\text{narr}}$.

EXAMPLE E. Assume that a h_λ -transformation of $(Y_i - \beta' x_i)/\sigma$ is $N\{0, 1\}$. When $h_\lambda(Z)$ is $N\{0, 1\}$, then Z has cumulative $\Phi(z)^\lambda$ and density $\lambda \Phi(z)^{\lambda-1} \phi(z)$. Hence Y_i has density

$$f(y_i, \sigma, \beta, \lambda | x_i) = \lambda \Phi((y_i - \beta' x_i)/\sigma)^{\lambda-1} \phi((y_i - \beta' x_i)/\sigma)/\sigma.$$

Is it now possible to evaluate partial derivatives w.r.t. σ, β, λ . Their null model versions, corresponding to $\lambda = 1$, become $(z_i^2 - 1)/\sigma, z_i x_i/\sigma, 1 + \log \Phi(z_i)$, where $z_i = (y_i - \beta' x_i)/\sigma$. Formula (3.4) gives the $(p+2) \times (p+2)$ matrix

$$J_{n,\text{wide}} = \begin{pmatrix} 2/\sigma^2 & 0 & b/\sigma \\ 0 & \sigma^{-2}n^{-1} \sum_{i=1}^n x_i x_i' & a\sigma^{-1}n^{-1} \sum_{i=1}^n x_i \\ b/\sigma & a\sigma^{-1}n^{-1} \sum_{i=1}^n x_i & 1 \end{pmatrix},$$

in which $a = EN \log \Phi(N) = 0.9032$ and $b = E\{1 + N^2 \log \Phi(N)\} = -0.5956$ (computed by numerical integration). It follows that

$$1/\kappa^2 = 1 - \frac{1}{2}b^2 - a^2 \bar{x}' \left\{ \frac{1}{n} \sum_{i=1}^n x_i x_i' \right\}^{-1} \bar{x}.$$

κ can be rather large, which in turn means that standard regression copes well even if λ differs quite a bit from 1. If only $|\lambda - 1| \leq \kappa/\sqrt{n}$, then standard regression methods work better than cumbersome ones employing a separate estimate for λ .

In the special case of a constant mean the tolerance limit against misspecification is very relaxed, with $\kappa = 12.090$. In this case $V = 1 + \log \Phi(z)$ is extremely well explained by $U = (z/\sigma, (z^2 - 1)/\sigma)$, with a maximal correlation of 0.993; see the discussion under 4F. The classic $N\{\xi, \sigma^2\}$ can stand a good deal of misspecification w.r.t. λ . — In another special case, that of $\alpha + \beta(x_i - \bar{x}) + \sigma Z_i$, $1/\kappa^2$ becomes $1 - \frac{1}{2}b^2$ and κ becomes much smaller, namely 1.103. In the language of 4F the n values of $V_i = 1 + \log \Phi(z_i)$ are now

much less well explained by the respective values of $U_i = (z_i, (x_i - \bar{x})z_i, z_i^2 - 1)/\sigma$, and the standard regression model can only tolerate up to $1.103/\sqrt{n}$ deviance from $\lambda = 1$.

EXAMPLE F. Write $p_i = p(x_i, \beta, \gamma)$, in which $\gamma = \gamma_0$ gives back ordinary logistic regression, and write p_i^0 for $p(x_i, \beta_0, \gamma_0)$ at some target point β_0 . It is not difficult to reach

$$J_{n,\text{wide}} = \frac{1}{n} \sum_{i=1}^n \frac{1}{p_i^0(1-p_i^0)} \begin{pmatrix} \partial p_i / \partial \beta \\ \partial p_i / \partial \gamma \end{pmatrix} \begin{pmatrix} \partial p_i / \partial \beta \\ \partial p_i / \partial \gamma \end{pmatrix}',$$

where the partial derivates are computed at the null point as usual. Finding the tolerance limit κ^2 is achieved by computing this $(p+1) \times (p+1)$ matrix numerically, at the target point, which will typically be the estimated $\hat{\beta}_{\text{narr}}$ computed from ordinary analysis, and then inverting it; κ^2 is found at the lower right corner.

In the two types of model departure discussed in Section 1, this goes as follows. If the wide model says $\alpha + \beta(x_i - \bar{x}) + \gamma(x_i - \bar{x})^2$, then

$$J_{n,\text{wide}} = \frac{1}{n} \sum_{i=1}^n p_i^0(1-p_i^0) \begin{pmatrix} 1 & t_i & t_i^2 \\ t_i & t_i^2 & t_i^3 \\ t_i^2 & t_i^3 & t_i^4 \end{pmatrix},$$

where $t_i = x_i - \bar{x}$. In the case of (1.5), on the other hand, involving a shape parameter η ,

$$J_{n,\text{wide}} = \frac{1}{n} \sum_{i=1}^n \frac{p_i^0}{1-p_i^0} \begin{pmatrix} (1-p_i^0)^2 & (1-p_i^0)^2 x_i & (1-p_i^0) \log p_i^0 \\ (1-p_i^0)^2 x_i & (1-p_i^0)^2 x_i^2 & (1-p_i^0) \log p_i^0 x_i \\ (1-p_i^0) \log p_i^0 & (1-p_i^0) \log p_i^0 x_i & (\log p_i^0)^2 \end{pmatrix}.$$

EXAMPLE G. Write $\sigma_1^2 = \sigma^2$ and $\sigma_2^2 = \sigma^2(1 + \gamma)$. Finding the J_{wide} matrix in the $(\xi_1, \xi_2, \sigma, \gamma)$ model is not difficult, and leads to $\kappa^2 = 2/\{r(1-r)\}$, where $r = m/(m+n)$. This means a tolerance level of $d = \kappa^2 J_{22} = 1/r = (m+n)/m$. The simple equal variance model can tolerate $\gamma^2 \leq 2(m+n)/mn$, which becomes $|\gamma| \leq 2/\sqrt{n}$ in the $m = n$ case. This is a fairly low tolerance limit, and different variances qualifies as a dangerous departure from the narrow model.

OTHER EXAMPLES. There is a large variety of other examples of common departures from standard models and that could be studied using our general methods and results. In each case one could compute the tolerance radius, one could speculate about robustness against the deviation in question in light of d and ρ of 4F, and one could implement the method of (5.3), for example. A partial list of such models and deviations is: (i) Typical i.i.d. models against various forms of dependence. (ii) Multinomial and log-linear models against higher order interactions. (iii) Analysis of variance models against interaction terms. (iv) Analysis of variance models against different variances in different groups. (v) Regression models against presence of cross-terms. (vi) Time series models against higher order autoregression or moving average. (vii) Typical i.i.d. models for discrete variables against Markov dependence. (viii) Markov chain models against second order Markovness. (ix) Models with normal errors against contamination of gross errors. (x) Traditional homogeneous models in survival analysis against heterogeneous frailness of individuals.

References

- Berger, J.O. (1982). Estimation in continuous exponential families: Bayesian estimation subject to risk restrictions and inadmissibility results. In *Statistical Decision Theory and Related Topics III*, eds. Berger and Gupta, 109–141. Academic Press, New York.
- Bickel, P.J. (1981). Minimax estimation of the mean of a normal distribution when the parameter space is restricted. *Ann. Statist.* **9**, 1301–1309.
- Bickel, P.J. (1983). Minimax estimation of the mean of a normal distribution subject to doing well at a point. *Recent Advances in Statistics, Festschrift for Herman Chernoff*, eds. Rizvi and Siegmund, 511–528. Academic Press, New York.
- Bickel, P.J. (1984). Parametric robustness: small biases can be worthwhile. *Ann. Statist.* **12**, 864–879.
- Efron, B. and Morris, C. (1971). Limiting the risk of Bayes and empirical Bayes estimators. *J. Amer. Statist. Assoc.* **66**, 807–815.
- Hjort, N.L. (1986a). *Theory of Statistical Symbol Recognition*. Research monograph, Norwegian Computing Centre, Oslo.
- Hjort, N.L. (1986b). Bayes estimators and asymptotic efficiency in parametric counting process models. *Scand. J. Statist.* **13**, 63–85.
- Hjort, N.L. (1991a). The exact amount of t-ness that the normal model can tolerate. Technical report, University of Oslo; submitted for publication.
- Hjort, N.L. (1991b). Computer programs and risk functions for estimators for a normal mean. Technical report, University of Oslo; available upon courteous request.
- Lehmann, E.L. (1983). *Theory of Point Estimation*. Wiley, Singapore.